

US007228441B2

(12) **United States Patent**  
**Fung**

(10) **Patent No.:** **US 7,228,441 B2**  
(45) **Date of Patent:** **\*Jun. 5, 2007**

(54) **MULTI-SERVER AND MULTI-CPU POWER MANAGEMENT SYSTEM AND METHOD**

(75) Inventor: **Henry T. Fung**, San Jose, CA (US)

(73) Assignee: **Huron IP LLC**, Grosse Pointe, MI (US)

(\* ) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 875 days.

This patent is subject to a terminal disclaimer.

4,398,192 A	8/1983	Moore et al. ....	340/825.44
4,463,440 A	7/1984	Nishiura et al. ....	364/900
4,479,191 A	10/1984	Nojima et al. ....	364/707
4,509,148 A	4/1985	Asano et al. ....	365/230
4,510,398 A	4/1985	Culp et al. ....	307/35
4,538,231 A	8/1985	Abe et al. ....	364/483
4,545,030 A	10/1985	Kitchin ....	364/900
4,570,219 A	2/1986	Shibukawa et al. ....	395/775
4,667,289 A	5/1987	Yoshida et al. ....	364/200
4,677,566 A *	6/1987	Whittaker et al. ....	700/295
4,698,748 A	10/1987	Juzswik et al. ....	364/200
4,766,567 A	8/1988	Kato ....	364/900

(Continued)

(21) Appl. No.: **09/860,210**

(22) Filed: **May 18, 2001**

(65) **Prior Publication Data**

US 2005/0177755 A1 Aug. 11, 2005

**Related U.S. Application Data**

(60) Provisional application No. 60/283,375, filed on Apr. 11, 2001, provisional application No. 60/236,043, filed on Sep. 27, 2000, provisional application No. 60/236,062, filed on Sep. 27, 2000.

(51) **Int. Cl.**  
**G06F 1/26** (2006.01)

(52) **U.S. Cl.** ..... **713/300; 713/320**

(58) **Field of Classification Search** ..... **713/300, 713/320**

See application file for complete search history.

(56) **References Cited**

**U.S. PATENT DOCUMENTS**

3,725,868 A	4/1973	Malmer, Jr. ....	395/375
4,279,020 A	7/1981	Christian et al. ....	364/900
4,316,247 A	2/1982	Iwamoto ....	364/200
4,317,180 A	2/1982	Lies ....	364/707
4,365,290 A	12/1982	Nelms et al. ....	364/200
4,381,552 A	4/1983	Nocilini et al. ....	364/900

Primary Examiner—James K. Trujillo

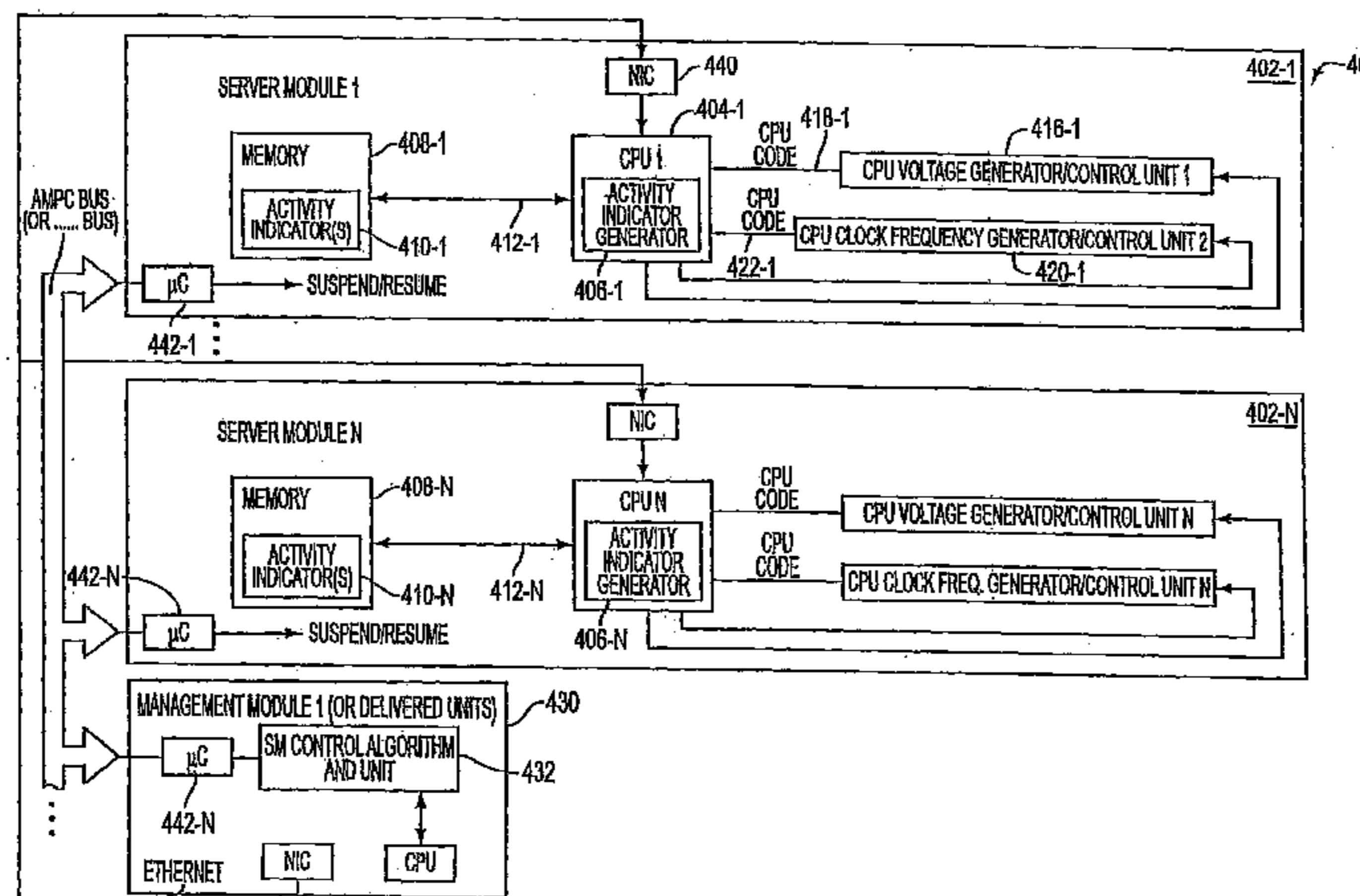
Assistant Examiner—Eric Chang

(74) Attorney, Agent, or Firm—Dykema Gossett PLLC

(57) **ABSTRACT**

Network architecture, computer system and/or server, circuit, device, apparatus, method, and computer program and control mechanism for managing power consumption and workload in computer system and data and information servers. Further provides power and energy consumption and workload management and control systems and architectures for high-density and modular multi-server computer systems that maintain performance while conserving energy and method for power management and workload management. Dynamic server power management and optional dynamic workload management for multi-server environments is provided by aspects of the invention. Modular network devices and integrated server system, including modular servers, management units, switches and switching fabrics, modular power supplies and modular fans and a special backplane architecture are provided as well as dynamically reconfigurable multi-purpose modules and servers. Backplane architecture, structure, and method that has no active components and separate power supply lines and protection to provide high reliability in server environment.

**35 Claims, 5 Drawing Sheets**



U.S. PATENT DOCUMENTS

4,780,843 A	10/1988	Tietjen .....	364/900	5,129,091 A	7/1992	Yorimoto et al. ....	395/750
4,809,163 A	2/1989	Hirosawa et al. ....	364/200	5,151,992 A	9/1992	Nagae .....	395/750
4,823,292 A	4/1989	Hillion .....	364/707	5,167,024 A	11/1992	Smith et al. ....	395/375
4,835,681 A	5/1989	Culley .....	364/200	5,175,845 A	12/1992	Little .....	395/550
4,841,440 A	6/1989	Yonezu et al. ....	364/200	5,201,059 A	4/1993	Nguyen .....	395/800
4,881,205 A	11/1989	Aihara .....	365/222	5,218,704 A	6/1993	Watts, Jr. et al. ....	395/750
4,893,271 A *	1/1990	Davis et al. ....	713/501	5,222,239 A	6/1993	Rosch .....	395/750
4,907,183 A	3/1990	Tanaka .....	364/707	5,247,164 A	9/1993	Takahashi .....	235/492
4,922,450 A	5/1990	Rose et al. ....	364/900	5,247,213 A	9/1993	Trinh et al. ....	307/465
4,963,769 A	10/1990	Hiltpold et al. ....	307/465	5,247,655 A	9/1993	Khan et al. ....	395/550
4,968,900 A	11/1990	Harvey et al. ....	307/296.3	5,249,298 A	9/1993	Bolan et al. ....	395/750
4,974,180 A	11/1990	Patton et al. ....	364/550	5,251,320 A	10/1993	Kuzawinski et al. ....	395/750
4,980,836 A	12/1990	Carter et al. ....	364/483	5,396,635 A	3/1995	Fung .....	395/800
4,991,129 A	2/1991	Swartz .....	364/707	5,710,929 A	1/1998	Fung	
4,996,706 A	2/1991	Cho .....	379/93	5,758,175 A	5/1998	Fung	
5,021,679 A	6/1991	Fairbanks et al. ....	307/66	5,799,198 A	8/1998	Fung	
5,025,387 A	6/1991	Frane .....	364/493	6,574,740 B1 *	6/2003	Odaohhara et al. ....	713/323
5,041,964 A	8/1991	Cole et al. ....	364/200	6,584,571 B1	6/2003	Fung	
5,083,266 A	1/1992	Watanabe .....	395/275	6,859,882 B2	2/2005	Fung	
5,119,377 A	6/1992	Cobb et al. ....	371/550	7,032,119 B2 *	4/2006	Fung .....	713/320
5,123,107 A	6/1992	Mensch, Jr. ....	395/800				

\* cited by examiner

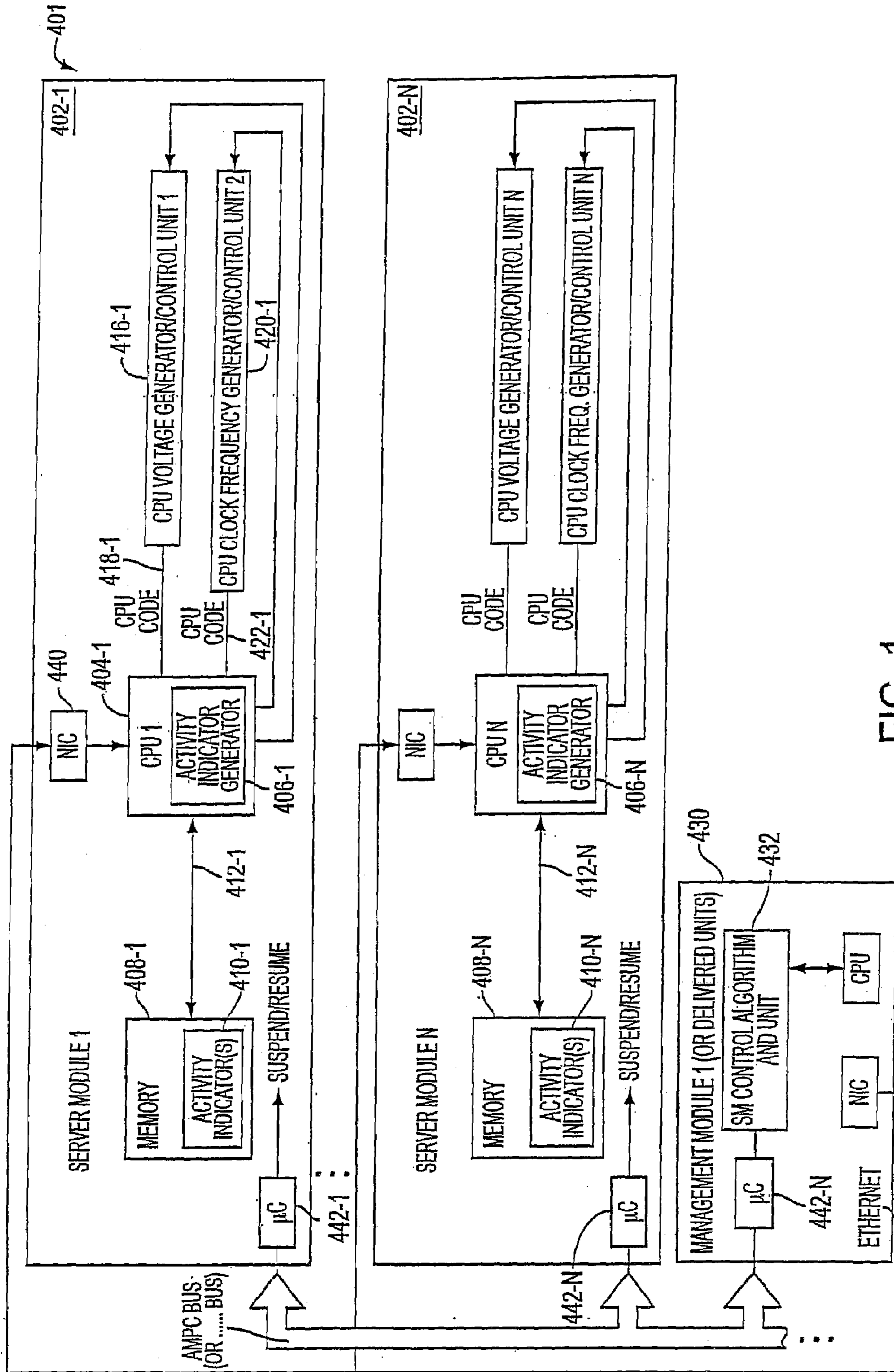


FIG. 1

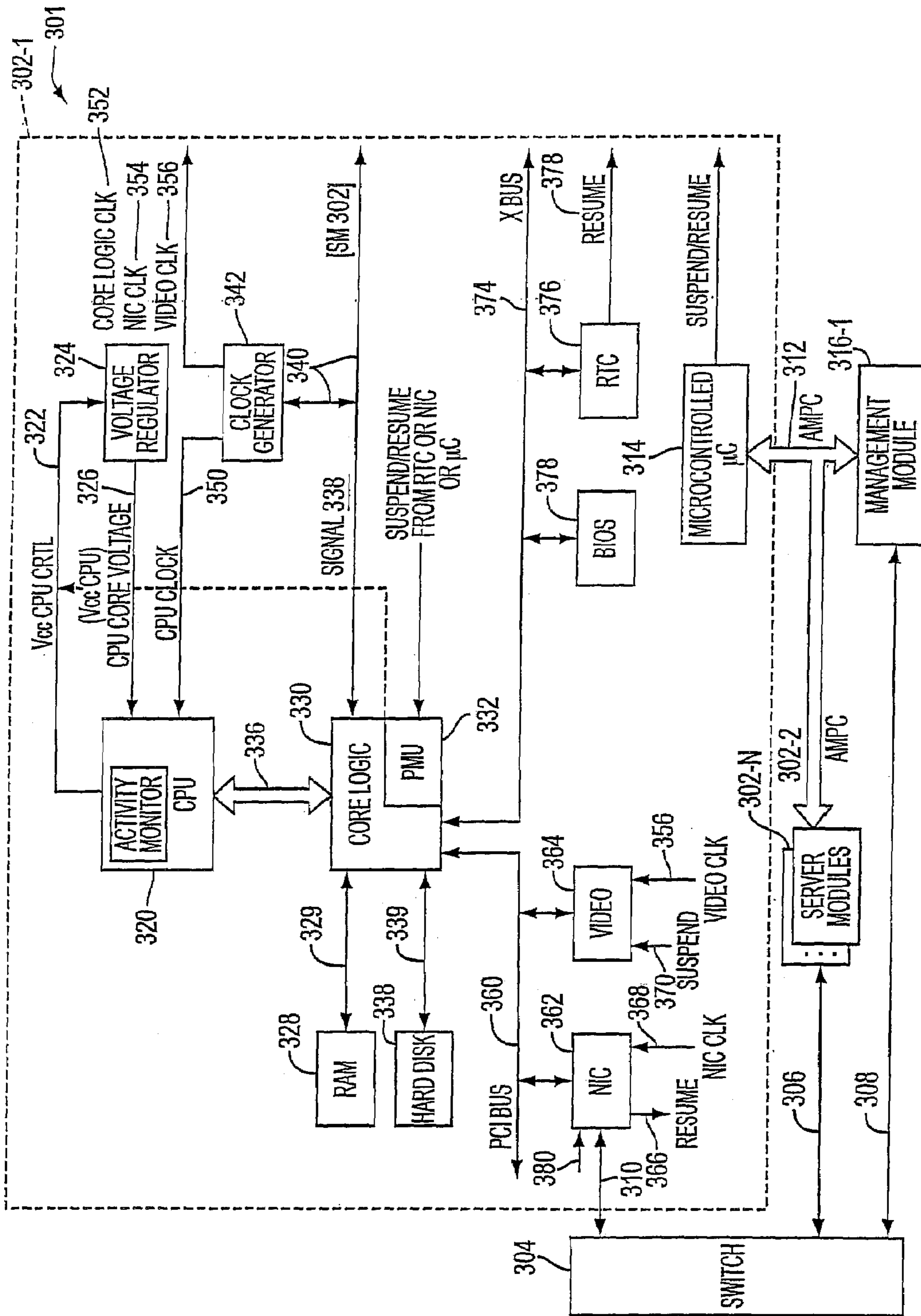


FIG. 2

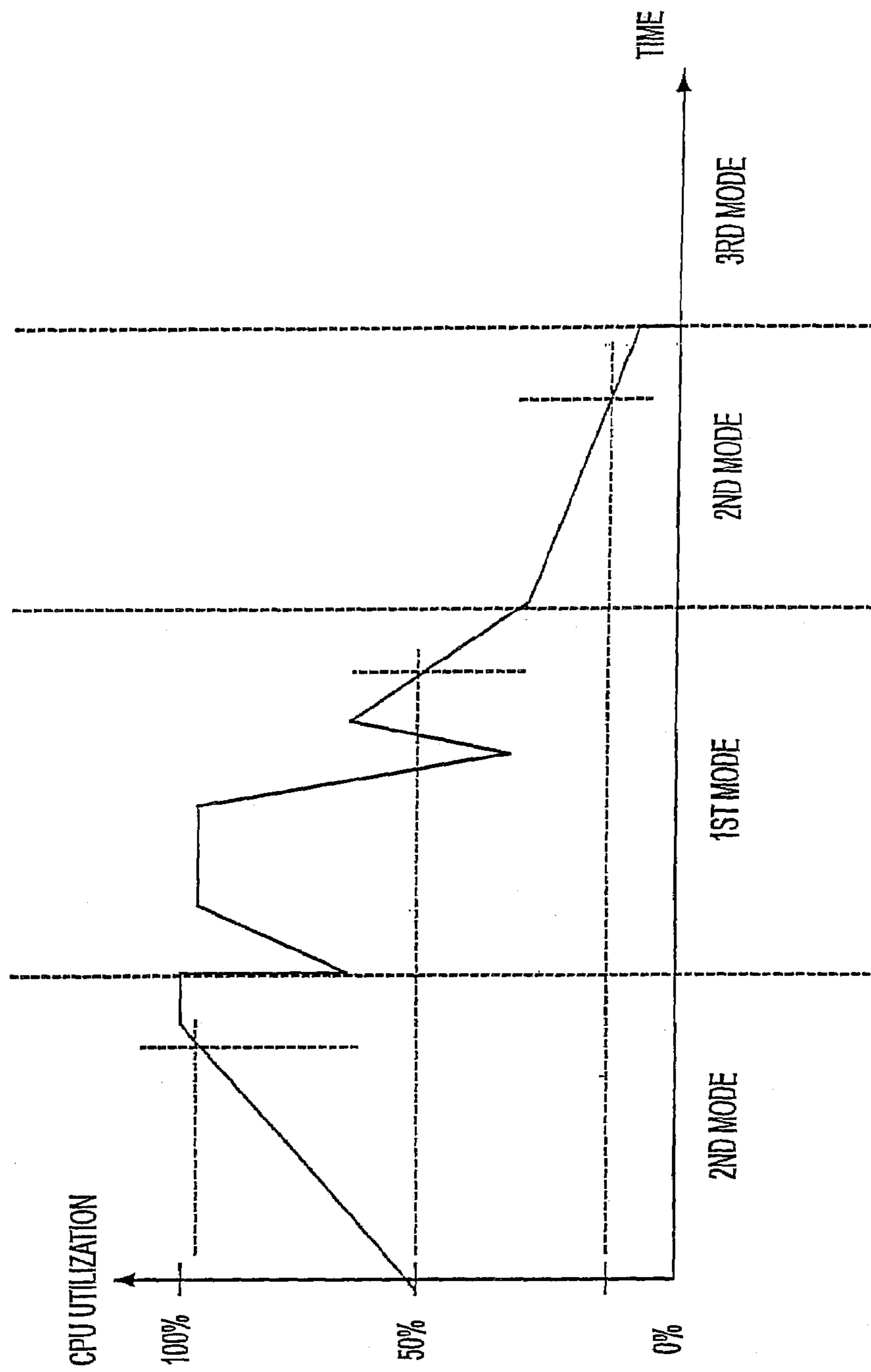


FIG. 3

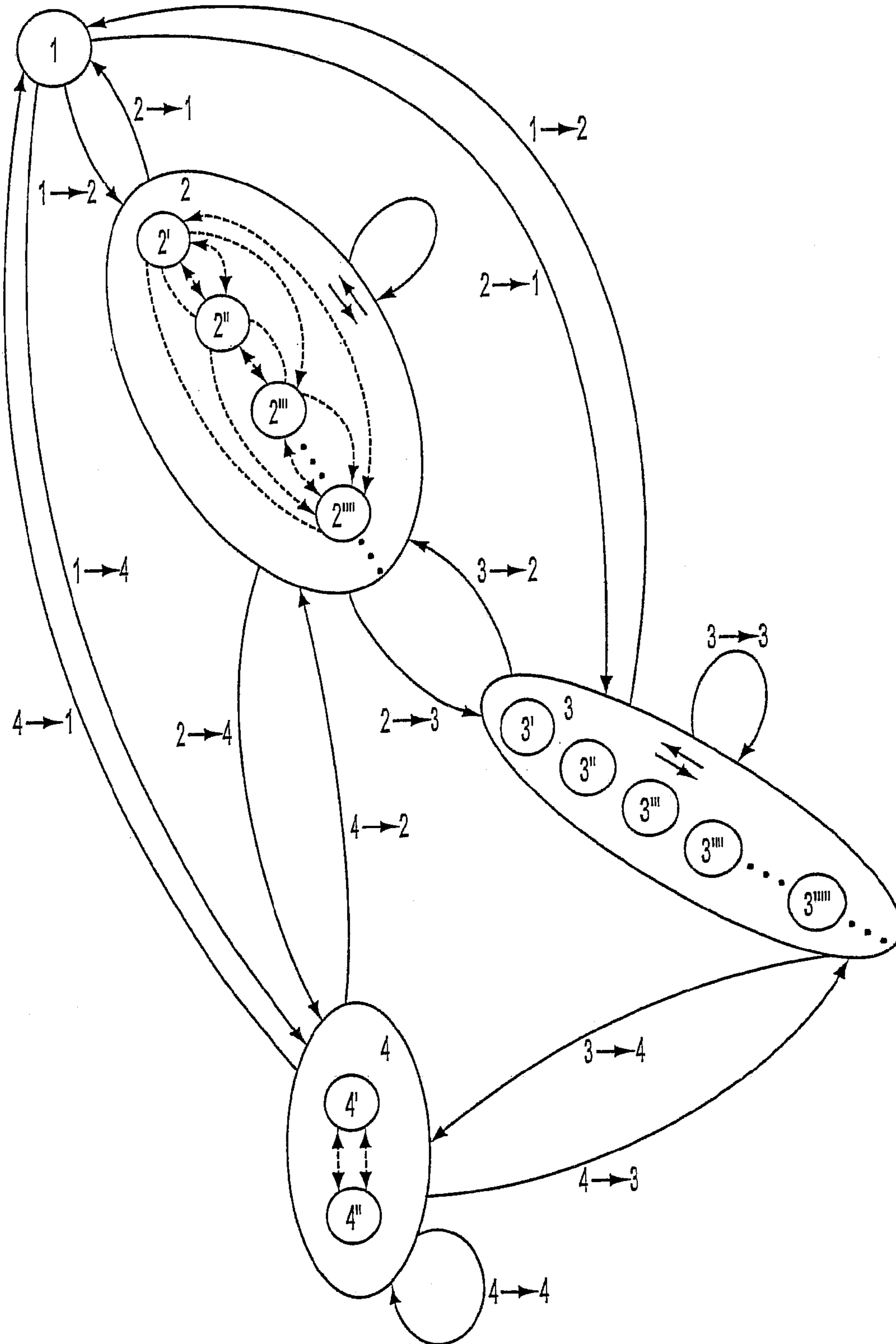


FIG. 4

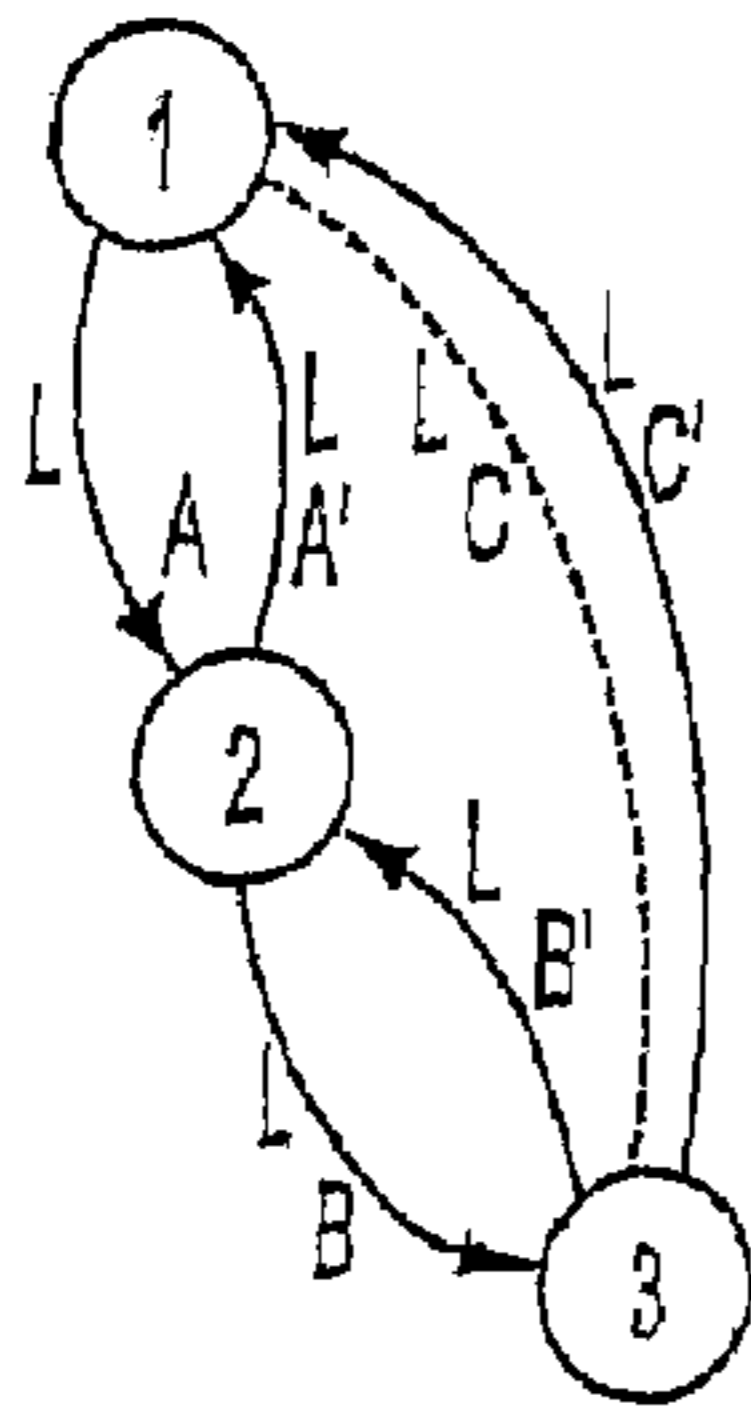


FIG. 5

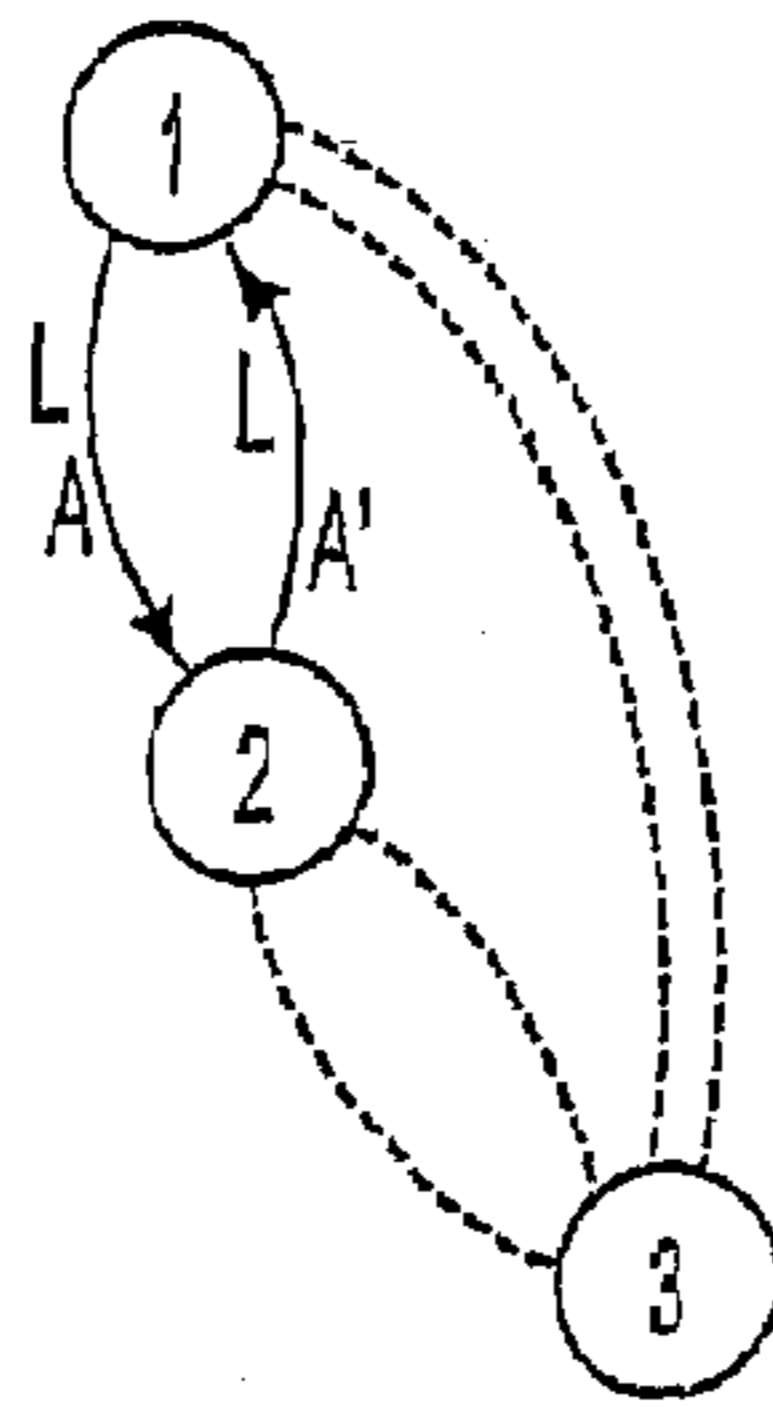


FIG. 6

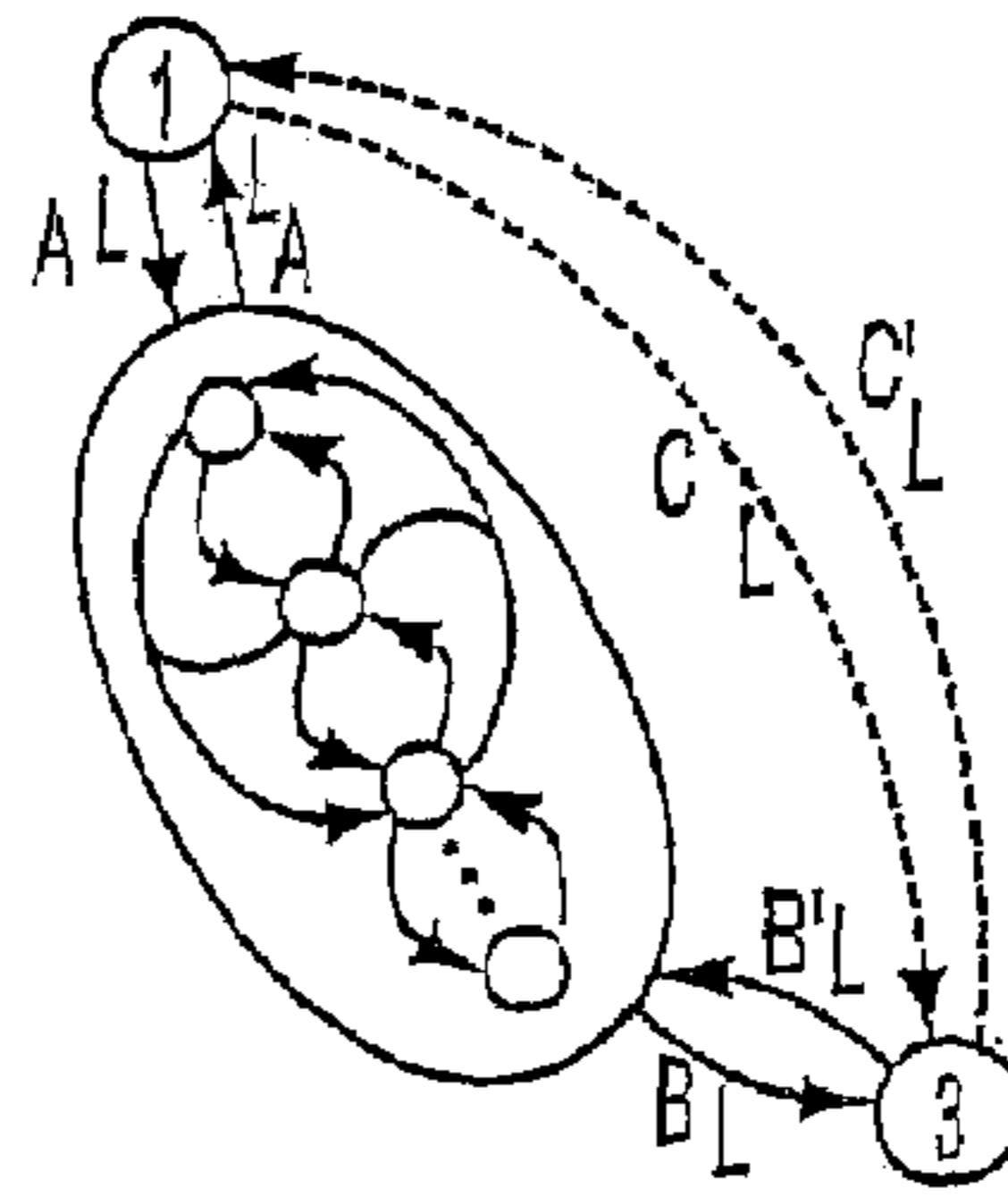


FIG. 7

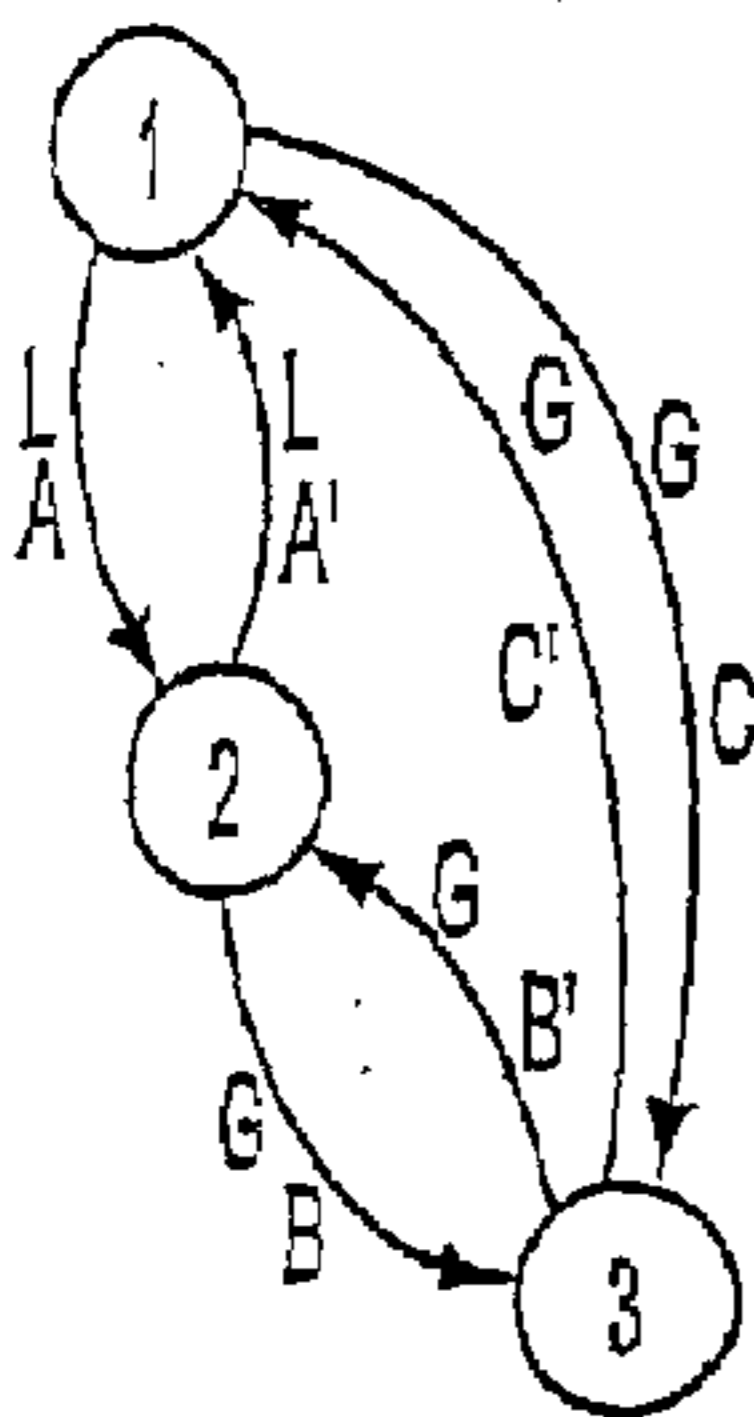


FIG. 8

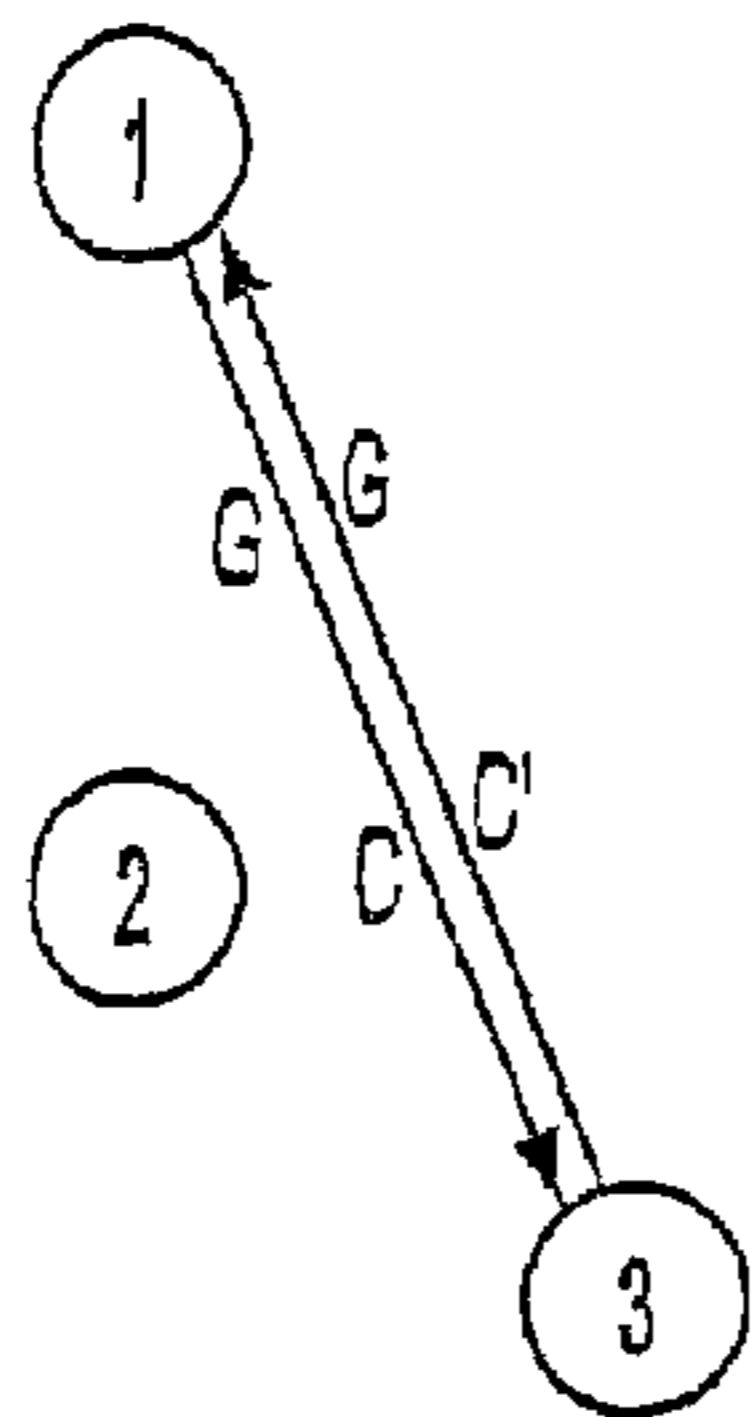


FIG. 9

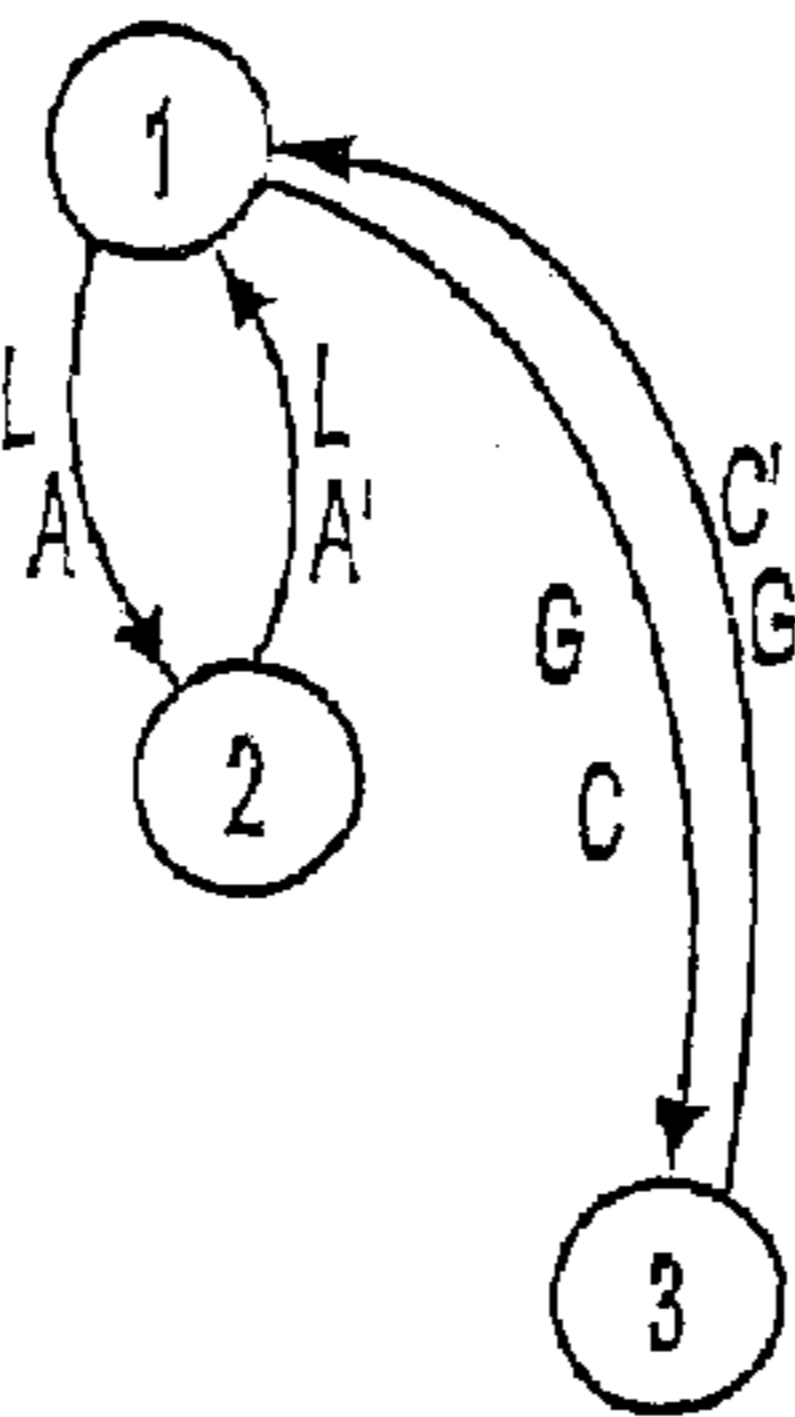


FIG. 10

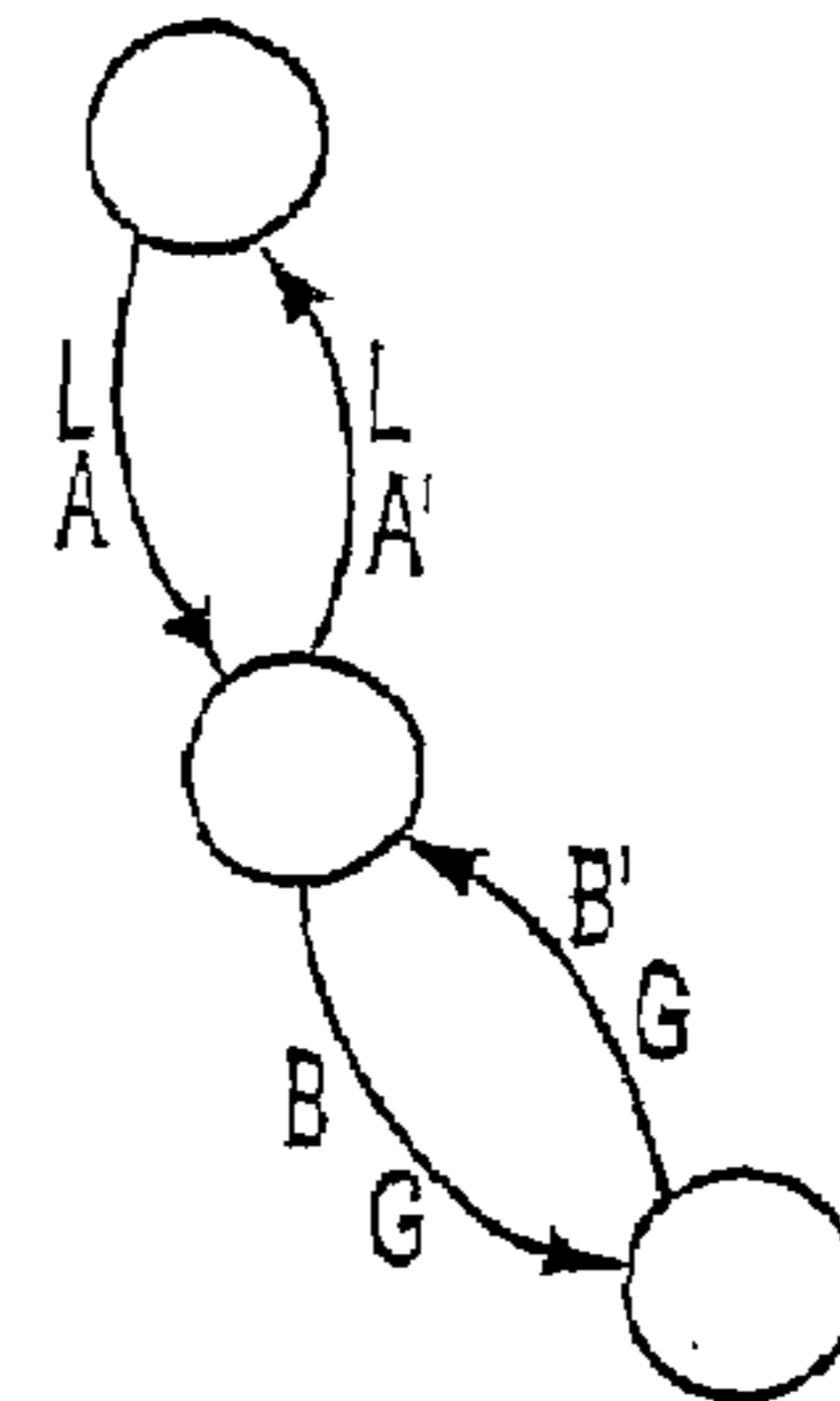


FIG. 11

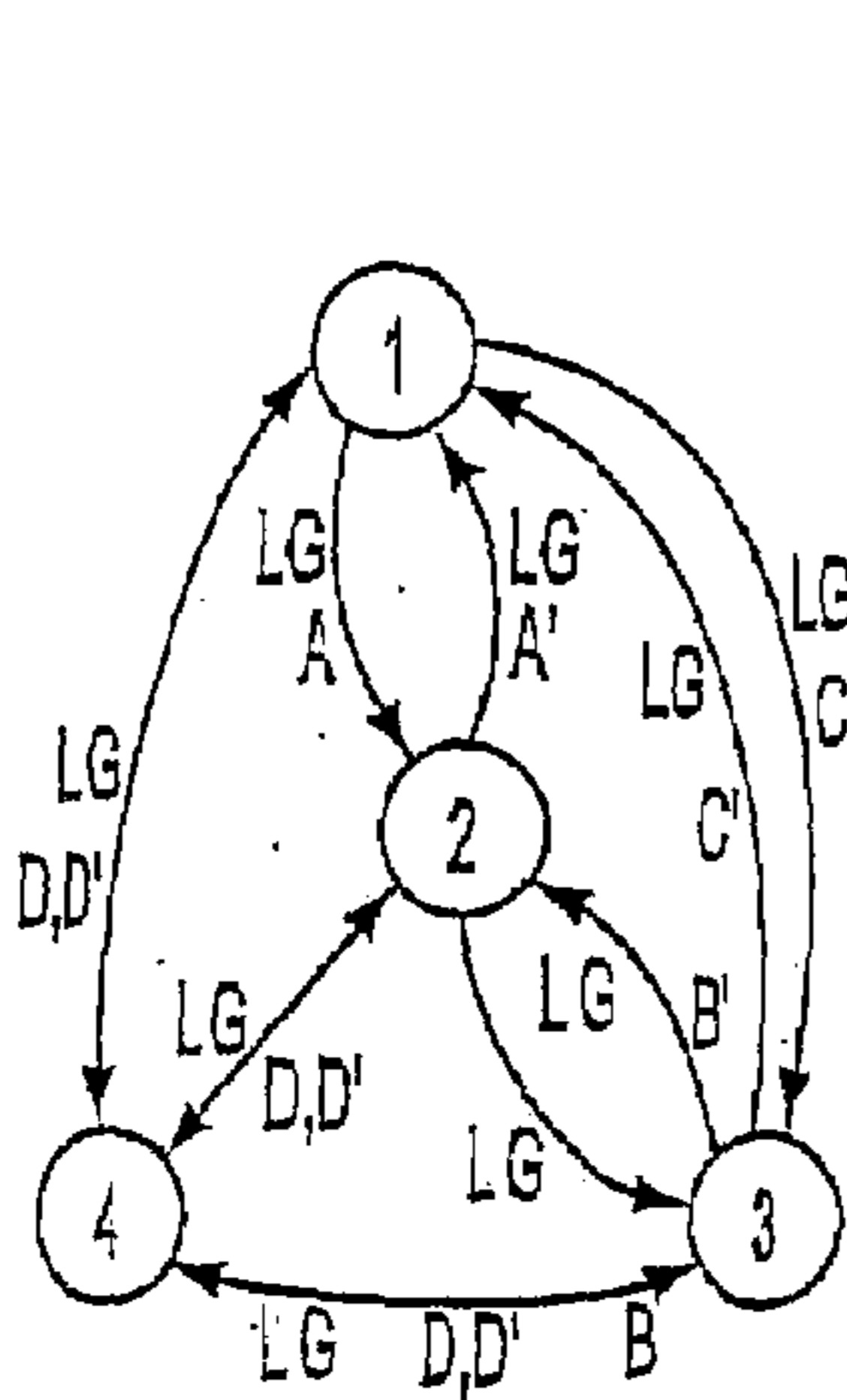


FIG. 12

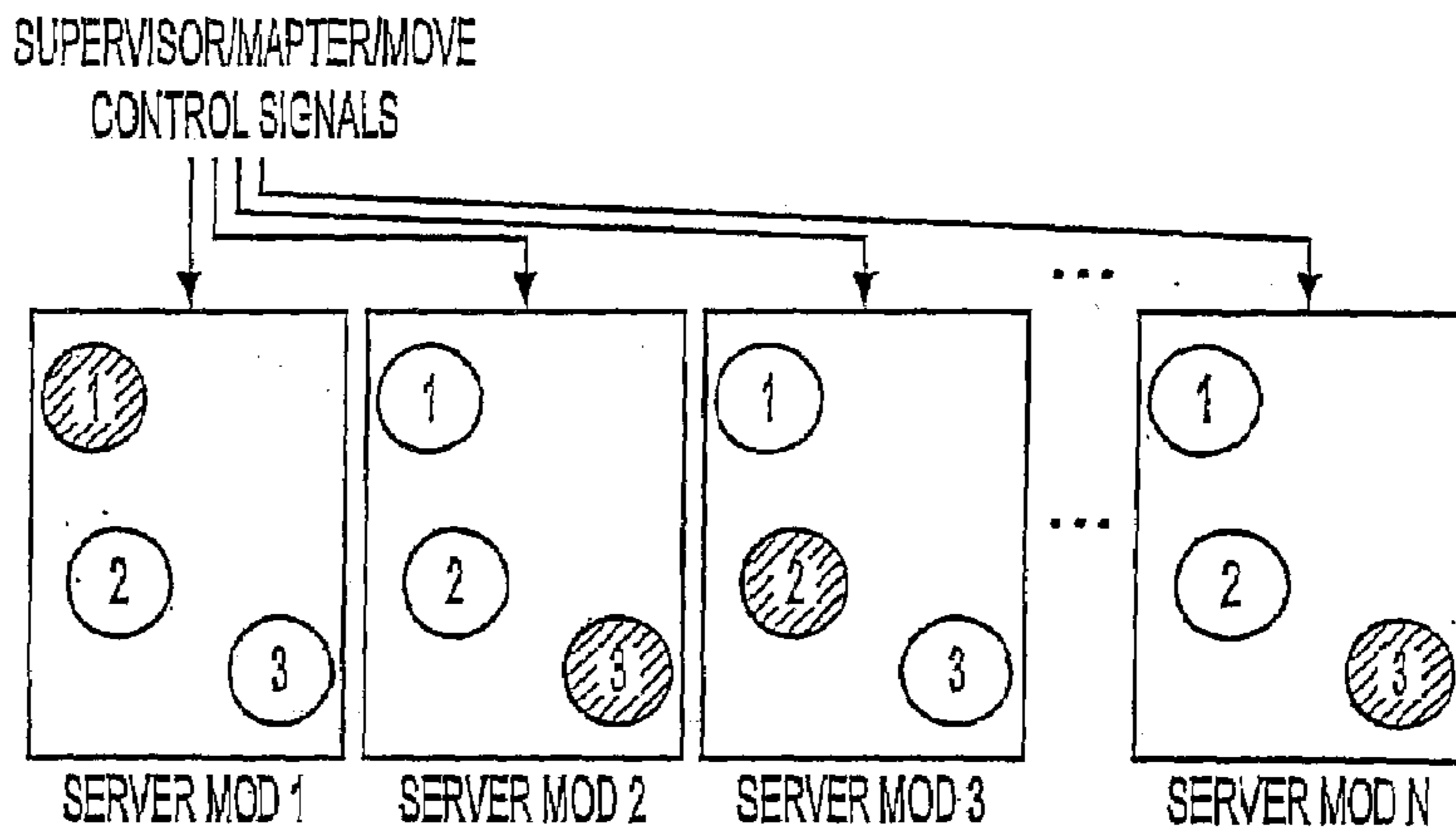


FIG. 13

## MULTI-SERVER AND MULTI-CPU POWER MANAGEMENT SYSTEM AND METHOD

### RELATED APPLICATIONS

This application is a continuing application under 35 U.S.C. §§ 119(e) and 120, wherein applicant and inventor claim the benefit of priority to U.S. Provisional Application Ser. No. 60/283,375 entitled System, Method And Architecture For Dynamic Server Power Management And Dynamic Workload Management for Multi-Server Environment filed 11 Apr. 2001; U.S. Provisional Application Ser. No. 60/236,043 entitled System, Apparatus, and Method for Power-Conserving Multi-Node Server Architecture filed 27 Sep. 2000; and U.S. Provisional Application Ser. No. 60/236,062 entitled System, Apparatus, and Method for Power Conserving and Disc-Drive Life Prolonging RAID Configuration filed 27 Sep. 2000; each of which application is hereby incorporated by reference.

### FIELD OF THE INVENTION

This invention pertains generally to architecture, apparatus, systems, methods, and computer programs and control mechanisms for managing power consumption and workload in data and information servers; more particularly to power consumption and workload management and control systems for high-density multi-server computer system architectures that maintain performance while conserving energy and to the method for power management and workload management used therein, and most particularly to system, method, architectures, and computer programs for dynamic server power management and dynamic workload management for multi-server environments.

### BACKGROUND

Heretofore, servers generally, and multi-node network servers in particular, have paid little if any attention to power or energy conservation. Such servers were designed and constructed to run at or near maximum levels so as to serve data or other content as fast as possible, or where service demands were less than capacity to remain ever vigilant to provide fast response to service requests. Increasing processor and memory speeds have typically been accompanied by higher processor core voltages to support the faster device switching times, and faster hard disk drives have typically lead to faster and more energy-hungry disk drive motors. Larger memories and caches have also lead to increased power consumption even for small single-node servers. Power conservation efforts have historically focused on the portable battery-powered notebook market where battery life is an important marketing and use characteristic. However, in the server area, little attention has been given to saving power, such servers usually not adopting or utilizing even the power conserving suspend, sleep, or hibernation states that are available with some Microsoft 95/98/2000, Linux, Unix, or other operating system based computers, personal computers, PDAs, or information appliances.

Multi-node servers present a particular energy consumption problem as they have conventionally be architected as a collection of large power hungry boxes interconnected by external interconnect cables. Little attention has been placed on the size or form factor of such network architectures, the expansability of such networks, or on the problems associated with large network configurations. Such conventional networks have also by-and-large paid little attention to the

large amounts of electrical power consumed by such configurations or in the savings possible. This has been due in part because of the rapid and unexpected expansion in the Internet and in servers connected with and serving to Internet clients. Internet service companies and entrepreneurs have been more interested in a short time to market and profit than on the effect on electrical power consumption and electrical power utilities; however, continuing design and operation without due regard to power consumption in this manner is problematic.

Networks servers have also by-and-large neglected to factor into the economics of running a network server system the physical plant cost associated with large rack mounted equipment carrying perhaps one network node per chassis. These physical plant and real estate costs also contribute to large operating costs.

In the past, more attention was given to the purchase price of equipment and little attention to the operating costs. It would be apparent to those making the calculation that operating costs may far exceed initial equipment purchase price, yet little attention has been paid to this fact. More recently, the power available in the California electrical market has been at crisis levels with available power reserves dropping below a few percent reserve and rolling blackouts occurring as electrical power requirements drop below available electrical power generation capacity. High technology companies in the heart of Silicon Valley cannot get enough electrical power to make or operate product, and server farms which consume vast quantities of electrical energy for the servers and for cooling equipment and facilities in which they are housed, have stated that they may relocate to areas with stable supplies of low-cost electricity.

Even were server manufactures motivated to adopt available power management techniques, such techniques represent only a partial solution. Conventional computer system power management tends to focus on power managing a single CPU, such as by monitoring certain restricted aspects of the single CPU operation and making a decision that the CPU should be run faster to provide greater performance or more slowly to reduce power consumption.

Heretofore, computer systems generally, and server systems having a plurality of servers where each server includes at least one processor or central processing unit (CPU) in particular have not been power managed to maintain performance and reduce power consumption. Even where a server system having more than one server component and CPU may possibly have utilized a conventional personal computer architecture that provided some measure of localized power management separately within each CPU, no global power management architecture or methods have conventionally been applied to power manage the set of servers and CPUs as a single entity.

The common practice of over-provisioning a server system so as to be able to meet peak demands has meant that during long periods of time, individual servers are consuming power and yet doing no useful work, or several servers are performing some tasks that could be performed by a single server at a fraction of the power consumption.

Operating a plurality of servers, including their CPU, hard disk drive, power supply, cooling fans, and any other circuits or peripherals that are associated with the server, at such minimal loading also unnecessarily shortens their service life. However, conventional server systems do not consider the longevity of their components. To the extent that certain of the CPUs, hard disk drives, power supplies, and cooling fans may be operated at lower power levels or for mechani-



cal systems (hard disk drive and cooling fans in particular) their effective service life may be extended.

Therefore there remains a need for a network architecture and network operating method that provides large capacity and multiple network nodes or servers in a small physical footprint and that is power conservative relative to server performance and power consumed by the server, as well as power conservative from the standpoint of power for server facility air conditioning. These and other problems are solved by the inventive system, apparatus and method. There also remains a need for server farms that are power managed in an organized global manner so that performance is maintained while reducing power consumption. There also remains a need to extend the effective lifetime of computer system components and servers so that the total cost of ownership is reduced.

### SUMMARY

Aspects of the invention provide network architecture, computer system and/or server, circuit, device, apparatus, method, and computer program and control mechanism for managing power consumption and workload in computer system and data and information servers. Other embodiments of the invention further provides power and energy consumption and workload management and control systems and architectures for high-density and modular multi-server computer systems that maintain performance while conserving energy and method for power management and workload management. Dynamic server power management and optional dynamic workload management for multi-server environments is provided by aspects of the invention. Modular network devices and integrated server system, including modular servers, management units, switches and switching fabrics, modular power supplies and modular fans and a special backplane architecture are provided as well as dynamically reconfigurable multi-purpose modules and servers.

### BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a diagrammatic illustration showing an embodiment of the inventive power conserving power management between two servers and a manager.

FIG. 2 is a diagrammatic illustration showing an alternative embodiment of a server system showing detail as to how activity may be detected and operating mode and power consumption controlled in response.

FIG. 3 is a diagrammatic illustration showing a graph of the CPU utilization (processor activity) as a function of time, wherein the CPU utilization is altered by entering different operating modes.

FIG. 4 is a diagrammatic illustration showing an exemplary state engine state diagram graphically illustrating the relationships amongst the modes and identifying some of the transitions between states or modes for operation of an embodiment of the inventive system and method.

FIGS. 5-12 are diagrammatic illustrations showing exemplary state diagram for operating mode transitions.

FIG. 13 is a diagrammatic illustration showing the manner in which a plurality of servers may operate in different modes based on local detection and control of selected mode transitions and local detection but global control of other selected mode transitions.

### DETAILED DESCRIPTION OF EMBODIMENTS OF THE INVENTION

The present invention pertains to computer system architectures and structures and methods for operating such computer system architectures in a compact high-performance low-power consumption manner. Computers, information appliances, data processing systems, and all manner of electronic systems and devices may utilize and benefit from the innovations described herein. Aspects of the invention also contribute to reliability, ease of maintenance, and longevity of the system as a whole and operation components thereof. In an application that is of particular importance and which benefits greatly from the innovations described here, the computer system is or includes a server system having at least one and more typically a plurality of servers. Each server will include at least one processor or CPU but may include multiple CPUs. In multiple server configurations significant power consumption reduction is achieved by applying the inventive power management scheme. These and other aspects of the invention are described in the sections that follow.

At least some embodiments of the invention provide a modular configuration where computers, servers, managers, and other devices and/or components are provided in a modular form so that such devices or components may readily be placed into service, maintained, removed from service, and/or configured within a rack or enclosure to provided desired operational features. References to "modular" devices, such as for example "modular server", "server module", "management module", or other module are therefore intended to apply to either a modular or non-modular device or component. For example, in the specification we conveniently refer to a "server" or "server module" to mean any server.

In still another aspect the inventive structure and method provide for significant power consumption reduction and energy savings as compared to conventional network and server architectures as only those power consuming resources that are actually needed to provide the quality of service required are in an active mode. Those node resources that are not needed may be powered off or placed in some power conserving standby mode until needed. In addition, operations performed by one or more nodes may be shifted to another node so that only the remaining active nodes consume power and the remaining nodes are in standby mode or powered off until needed. The intelligence within one of the nodes acting as a master node for the cluster or ISS may then wake up the inactive node and configure it for operation. A system may be woken up and placed in any of the available operating modes by any one of a plurality of events. Nodes may also be placed into an inactive or power conserving mode when no demands are made on their resources independent of whether responsibility for their functionality has been shifted to another node or nodes. In one embodiment of the invention the power consumed is reduced by a factor of about 10-times as compared to a standard 19-inch wide by 1.75-inch high (1U) rack mountable network node device. This power savings is accomplished at least in part by one or more of the following measures: the reduction in the number of power supplied, use of the mounting plate as a heat sink to assist in removing heat from the enclosure, providing power saving controls to circuits and devices within the ISS enclosure, and the above described ability to reconfigure and take off line unneeded capacity.

## 5

Many different types of servers architectures are known in the art. Typically, such servers have at least one processor with associated fast random access memory (RAM), a mass storage device that stores the data or content to be served by the server, a power supply that receives electrical power (current and voltage) from either a battery or line voltage from an electrical utility, a network communication card or circuit for communicating the data to the outside world, and various other circuits that support the operation of the CPU, such as a memory (typically non-volatile ROM) storing a Basic Input-Output System (BIOS), a Real-Time Clock (RTC) circuit, voltage regulators to generate and maintain the required voltages in conjunction with the power supply, and core logic as well as optional micro-controller(s) that communicate with the CPU and with the external world to participate in the control and operation of the server. This core logic is sometimes referred to as the Northbridge and Southbridge circuits or chipsets.

From a somewhat different perspective, variations in server architecture, reflect the variations in personal computers, mainframes, and computing systems generally. The vast structural, architectural, methodological, and procedural variations inherent in computer systems having chips, chipsets, and motherboards adapted for use by Intel Processors (such as the Intel x86, Intel Pentium™, Intel Pentium™ II, Intel Pentium™ III, Intel Pentium™ IV), Transmeta Crusoe™ with LongRun™, AMD, Motorola, and others, precludes a detailed description of the manner in which the inventive structure and method will be applied in each situation. Those having ordinary skill will appreciate in light of the description that the inventive structure and method apply to a broad set of different processor and computer/server architecture types and that minor variations within the ordinary skill of a practitioner in the field may be made to adapt the invention to other processor/system environments.

Before describing particular implementations that relate to more or less specific CPU designs and interfaces, attention first directed to a simplified embodiment of the inventive system and method with respect to FIG. 1. In this embodiment, at least two (and up to n) servers or server modules (where servers are made in modular form or configuration) 402-1, . . . , 402-N are provided, each including a CPU 404 and a memory 408. CPU 404 includes an activity indicator generator 406 which generates activity indicators, and either (i) communicates the activity indicators to memory 408 for storage in an activity indicator(s) data structure 410, or not shown, (ii) communicates them directly to a server module control unit and algorithm 432 within management module 430. Different types of activity indicators such as are described elsewhere in the specification, such as for example an idle thread based activity indicator may be used. Whether stored in memory or communicated directly, the activity indicator(s) are used by the management module to determine the loading on each of the server modules individually and as a group. In one embodiment, activity information or indicators created on any one computer or device (such as a server module) is accessible to a manager or supervisor via standard networking protocol.

Although not illustrated in FIG. 1, analogous structure and signals generated and received may be used to control the operation of core logic circuits to thereby control core logic voltage and core logic clock signals in a manner to reduce power consumption where such core logic power management is provided.

## 6

Voltage and frequency are regulated locally by the CPU using an activity monitoring scheme, such as for example one of the activity monitoring scheme illustrated in Table I.

TABLE I

Exemplary Activity Monitoring Schemes carried out in CPU or PMU

	Carried out by CPU	Carried out by PMU
Application Layer	Port Address	NA
Network Layer	TCP/IP	NA
Physical Layer	Idle Threads, Activity Counter	I/O Activities

This power management scheme may be interpreted in one aspect as providing a Mode1-to-Mode2 and Mode2-to-Mode1 power management scheme, where both Mode 1 and Mode 2 are active modes and the state of the CPU in either Mode 1 or Mode 2 is controlled locally by the CPU, and in another aspect as providing a Mode3 (inactive mode or maintenance of memory contents only). Mode3 control may also be performed locally by the CPU, but in one of the preferred embodiments of the invention, entry into a Mode 3 stage is desirably controlled globally in a multi-CPU system. Where the multi-CPU's are operative with a plurality of servers for multi-server power management, the Management Module (or a Server Module acting as a manager on behalf of a plurality of server modules) determines which Server Module should enter a Mode 3 state using the Server Module control algorithm and unit 432. Activity monitoring of individual Server Modules 402 is desirably based on the standard network protocol, such as for example SNMP. Therefore the activity indicators may be retrieved from the CPU 406 or memory 408 via NIC 440 as is known in the art. A communication link coupling micro-controllers (1C) 442 together, and in particular the micro-controller of the Management Module with the microcontrollers of the several Server Modules. This permits the management module to communicate commands or signals to the server modules which are received by the microcontrollers even when the CPUs are in a suspended state (Mode 3). In so providing for monitoring over the first link (the Ethernet) and control over the second link (the AMPC bus), the server modules may be monitored for activity and controlled globally to reduce power consumption while providing sufficient on-line capacity. It is noted that the power management may be effected by altering either or both of the CPU clock frequency 420 or the CPU voltage 416.

Although a separate management module 430 is illustrated in FIG. 1, it should be understood that the management functionality generally, and the server module control algorithm in particular may be implemented by one of the operating server modules. For example, the control algorithm would be implemented as a software or firmware procedure executing in the CPU and processor of a server module designated according to predetermined rules, policies, or procedures to be the master.

It is noted that although several of the modes described conserve power, they do not compromise performance, as the cumulative combination of server modules is always maintained at or above minimum targeted performance.

In FIG. 2 there is illustrated an exemplary system 301 including a server (such as for example, an ISSU server module) 302-1, coupled to a switch (such as for example, an ISSU switch module) 304, and through the switch 304 and optionally via a micro-controller ( $\mu$ C) 314 within server 302 over a separate (optional) direct bus connection 312 (such as

for example, the AMPC bus made by Amplus of San Jose, Calif.) to a power management supervisor (such as for example, ISSU management module) **316**. As described elsewhere herein, switch **304** is responsible for connecting the various server module(s) **302**, management module(s) **316**, and other components that are or may be controlled to achieve the power conservation features of the invention. Recall that such subsystems as the power supply (not shown) and cooling or fan modules may also be coupled through the switch **304**. The connectivity and signals shown in the diagram are intended to show significant control paths pertinent to the operation of the invention, and therefore some signals that are conventional or do not illustrate the operation of the invention are not shown to avoid obscuration of the invention.

Attention is now focused on the internal structure and operation of the server module **302**. During operation CPU **320** executes commands or instructions, or when no instructions are present to be executed, executes idle threads. The activity level of the CPU is monitored and a control signal **Vcc\_CPU\_control 322** is generated based on that sensed activity or lack of activity. The manner in which this activity is sensed or the manner and characteristics of the **Vcc\_CPU\_control** signal will typically vary depending on the processor type, operating system, and other factors specific to the system architecture. By way of illustrative example, an indication as to the CPU activity or lack of activity may be generated by monitoring by executing an application layer function call that returns a value indicating the idle thread execution based activity. This is possible in the Microsoft Windows 98, 2000, and NT operating environments, for example.

As the name implies, the **Vcc\_CPU\_control** signal **322** which is an input signal to voltage regulator **324** controls or influences the CPU core voltage **Vcc\_CPU 326**. As described elsewhere in this description, the CPU core voltage **326** may be raised and lowered in conjunction with the CPU clock frequency to provide adequate switching response of CPU circuits without excessive voltage. Although this embodiment illustrates that the **Vcc\_CPU\_control** signal **322** is generated within the CPU, in an alternative embodiment, it may be generated within the core logic block **330**. In one embodiment, the CPU clock is adjusted based on a signal from the core logic and the CPU voltage is adjusted on the basis of the CPU itself. This is due to the fact that the voltage change is desirably synchronized in time with the frequency change. In some sense, this control may be viewed as including an effective link from the core logic to control the voltage regulator output.

Core logic **330** includes a Power Management Unit **332** of which many types are now known; however, one early example of a Power Management Unit is described in U.S. Pat. Nos. 5,396,635, 5,892,959 and 6,079,025 (each of which is herein incorporated by reference) by the inventor of the present application as well as in the other applications related thereto. In operation, PMU **332** receives a signal over bus **336** and generates an output signal **338** that is communicated overbus **340** to clock generator **342**. Clock generator block **342** includes circuitry that generates a CPU clock **50**, a core logic clock signal **352**, a Network Interconnect Card (NIC) clock signal **354**, and a video clock signal **356**.

RAM **328** is coupled to core logic **330** via DRAM control line and hence to the CPU via bus **336**. Hard disk drive **338** is similarly coupled to core logic **330** to CPU via bus **336**. In one implementation, Redundant Array of Independent Disc (RAID) data storage is provided for the server modules. As is known, this RAID storage provides considerable data

redundancy. In order to implement this RAID in a power management efficient manner, two IDE controllers (or enhanced IDE controllers) are used to interface to two separate disk drives. Provision of two hard disk drives supports RAID Level 0, RAID Level 1, and RAID Level 0+1 implementations. Aspect of the RAID power management disk drive longevity are described in co-pending U.S. Provisional Application Ser. No. 60/236,062 entitled System, Apparatus, and Method for Power Conserving and Disc-Drive Life Prolonging RAID configuration filed 27 Sep. 2000, hereby incorporated by reference. It is noted that providing RAID storage or multiple disk drives on the servers is advantages though not required.

Clock generator **342** includes clock signal generating and logic circuitry or other means for generating a CPU clock signal at the desired frequency or for selecting a CPU clock signal from an available plurality of clock signal having different frequencies. Under the inventive power management scheme, the clock frequency is adjusted downward within a permissible CPU clock frequency range to provide a CPU processing power that matches the present need, and to the extent that the present need is below the maximum capability of the processor when operating at full permissible clock frequency, to reduce the power consumption of the CPU. As the CPU core voltage may be reduced below a maximum voltage when the clock frequency is below its maximum frequency, the CPU core voltage may be lowered with the clock frequency or speed.

A PCI bus **360** coupling NIC **362** and Video processor **364** is provided and interfaces with CPU **320** via Core logic **330**. NIC **362** generates and provides a resume output **366** and NIC Clock input signal **368**, and Video processor **364** is provided with a video clock signal **356** from the clock generator **342** and a suspend input signal **370**. It is noted that the suspend and resume signals may come from multiple sources to affect the desired control and management.

In this illustrative embodiment, an X-bus **374** is provided to couple the Real-Time Clock (RTC) **376** and BIOS **378** to the core logic **330** and via bus **336** to the CPU as required. RTC **376** may generate a resume output signal **378**. This RTC generated resume signal **378** is therefore operative to activate PMU **332**, core logic **330**, and CPU **330** under a predetermined time or alarm condition. For example, the RTC may be set to generate a resume signal **378** at 8:00 am local time every day to bring the server module **302** back online.

The NIC resume signal may be generated when a specific packet is received. When generated in one of these manners and communicated to the PMU **332** it is operative to place the core logic **336** back into an active state and hence CPU **320** into any selected state of mode. One situation in which the NIC resume signal may be generated is when the server module is in a powered-on but inactive state, such that the CPU clock is stopped (or operating at an extremely low clock frequency). Under such condition, a simple way of waking the server module **302** is to communicate a signal **380** from management module **316** via switch **304**. As the NIC will typically be kept active, it will receive the signal **380** and generate the resume signal **366**.

It is noted that each of the elements, such as the hard disk drive, Video processor and other power consuming elements may include means for receiving a control signal that places them into a power conserving state or that brings them out of on or more power conserving states into a full power and performance mode.

It is noted that the embodiment illustrated in FIG. 2 represents a system that might utilize any of a number of

conventional processors or CPU, and might for example utilize a CPU of the Intel Pentium, Pentium II, Pentium III, or Pentium IV types made by Intel Corporation of Santa Clara, Calif., various Advanced Micro Device CPUs, CPUs made by Transmeta, as well as other processors and CPUs as are known in the art.

Having now described the physical architecture and connectivity of an exemplary Integrated Server System, the structure and operation of an exemplary server module, management module, and switch module, aspects of how these modules work independently and in synergistic manner to accomplish significant power or energy conservation without sacrificing performance (or with only an undetectable insignificant performance impact at most) are now described.

Conventional servers do not take power consumption or power savings into consideration in their normal operation. The primary philosophy of data center and internet service providers (ISPs) is over provision. If one considers the relationship between traffic (or load) and the power consumption, conventional servers operate at a relatively constant and high power level that is independent of load. The processors and associated memory typically run at maximum clock rate or frequency, full supply voltage to the processor core, hard disc drives on and rotating constantly, and video and other sub-systems and ports on and operable all the time independent of whether they are being used at that time.

Power conservation features such as may be provided in consumer personal computers (PCs) such as the suspend, sleep, hibernation, and the like types of reduced power operation. Several different power management interface specifications have been developed in recent years, including for example, the Advanced Configuration and Power Interface Version 1.0 (herein incorporated by reference) co-developed by Intel, Microsoft and Toshiba, which specifies how software and hardware components like the operating system, motherboard and peripheral devices (such as hard disk drive) talk to each other about power usage.

One embodiment of the inventive structure is directed as so called "front end server" applications. While the power consumption of conventional servers may vary depending upon the processor type and operating characteristics, number of processors, amount of memory (RAM), disc drive or other storage device type and number, and the like, most conventional servers such as those made by Cobalt, Compaq, Dell, and others consume some where in the range of between about 50 Watts to 150 Watts or more. Some servers have as many as four processors and will consume considerably power.

Conventional servers typically have similar architecture to personal computers made for home and business use, albeit with greater I/O capabilities and horsepower for their intended tasks. Interestingly, most of these servers retain the video capabilities in spite of the fact that the servers will not be used by anyone for viewing the video. It is unfortunate that the video circuitry (either within the processor or as a separate chip) is consuming power yet produces no beneficial effect or result.

The structures and methods of the invention provides a very low power design so that even when the inventive server is operating at its maximum performance level and consuming its maximum power, that maximum power consumption is still a fraction of the maximum (and steady-state) power consumption of conventional non-power managed processors and servers. This maximum power level is typically between about 10 to 15 Watts though it may fall

within other ranges or be reduced further. This reduction is possible for several reasons, including the provision of a very low power consumption processor or CPU, turning off devices or components within the system that are not being used at the time. Another significant power savings is provided by power managing the CPU according to the network traffic or server load conditions. Therefore the power consumption is less than the maximum power consumption unless the load is at a peak and all of the devices and components are powered on to handle the load. With this throttling back as a function of load, the power consumption may be at any intermediate value between zero (when and if the unit is powered off completely) or at a very low power consumption level when placed in some power conserving mode (such as a sleep, suspend, or other specialized power conserving mode as described elsewhere herein). Thus, capabilities of the server are matched to the demands being placed on the server. This power control or management is referred to as power on demand (Power on Demand™) and permits power conservation without any loss of server capability. Power management may also be controlled dynamically.

The over-provisioning of servers by ISPs and Data Centers is adapted at least in part because e-commerce can be highly seasonal and subject to considerable event driven demand surges. For example, the traffic or load requirements placed on servers during Christmas Holiday season may be many time or even one or more orders of magnitude as compared to other times of the year. News, stock market, and other organizations may have analogous traffic fluctuations during a single day. Unless such e-commerce entities are able to satisfy the inquiries of their customers with tolerable quality of service (QOS), such customers may never come back to the site. Therefore, day-to-day, week-to-week, and month-to-month traffic loading can vary over a wide range. For one typical ISP, the average load is about twenty-percent (20%) of the maximum load.

In the inventive system and method, by varying the power consumption according to load, considerable additional savings are realized. For an exemplary system in which the base maximum power consumption is 10 watts rather than 50 watts, and the power consumed during the service cycle is on average 20% of the maximum, the net result is a realization of the product of these two savings for a savings of about 25 times. That is the power consumed over a day is  $\frac{1}{25}$  of the power consumed for a conventional server operation.

Typically, the amount of power savings and then relationship between traffic and power consumed will depend upon the nature of the server. For example, a web server may exhibit a different load versus power consumption characteristic curve than a streaming video server, which will be different that a content caching server. These relationships may be linear or non-linear. The nature of the content may also impact this relationship.

The inventive scheme interactively reacts to the load and scales the number of components and/or devices as well as the operational parameters and operating characteristics of the devices and/or components to match the load or a predetermined quality of service, or some other identified performance target for the server system.

The inventive ISS may incorporate multiple servers adapted to serve different type of content. Thus it may be expected that each different server will exhibit somewhat different power consumption and power consumption reduction characteristics. These characteristics need not be known a priori to realize their benefits.

Attention is now directed toward a description of exemplary different operating modes. In one aspect the inventive structure and method provide for a transition in a single processor or CPU between a first mode (Mode 1) and a second mode (Mode 2) wherein the second mode consumes less power or energy than the first mode. Power or energy consumption in the processor or CPU (and optionally in other circuit components or peripherals connected to or associated with the processor or CPU) may be reduced in a variety of ways, including for example, lowering a processor or CPU core voltage, reducing a processor or CPU clock frequency, or lowering the core voltage and the clock frequency at the same time.

In some systems and methods, the core voltage and clock frequency are changed continuously or in stages in some synchronized manner, as a higher core voltage may typically be required to support a faster processor or CPU clock frequency. It is noted that the first and second mode are each active operating modes in which the processor or CPU is executing instructions and carrying out normal processor functions. While the core voltage may be reduced, the processor clock is still cycling at some nominal rate. The lower limit for processor clock frequency reduction may generally be selected based on the types of processing that may be accomplished at that rate. For example, first mode operation would typically be provided at substantially 100% of the nominal rated clock frequency for the processor, while second mode operation provide a clock frequency less than 100%. Such reduced processor clock frequency may generally be in the range of between about 5% to 95% of the maximum, more usually between about 20% and about 80%, more usually between about 20% and 60%. In some systems, the processor clock may be reduced by factors of two using clock signal division circuitry. In other systems, the processor clock frequency may be reduced in fixed increments or according to a clock frequency rate reduction lookup table or algorithm in a clock generator circuit. As the second mode may be considered to be any active operating mode less than the first mode, it will be understood that there may be multiple levels of this second mode. That is, Mode 2 may be multi-level.

In addition to these first and second modes, the processor or CPU may be placed into an inactive third mode (Mode 3) characterized by consuming less power or energy (conserving more power or energy) than in the first mode or the second mode. This third mode is referred to as an inactive mode as the processor clock will be stopped or operate at such a low frequency that the processor effectively processes no instructions and performs substantially no useful work relative to the amount of work provided in the first or second modes. Usually, the processor clock will be stopped and where core voltage control is available, the processor core voltage will be reduced to a level just sufficient to maintain processor state. This third mode is distinguished from a fourth mode (Mode 4) where the processor is powered off and does not maintain processor state, revitalization of the processor from the fourth mode requiring a reboot or other initialization procedure. Such reboot or initialization procedures typically requiring a few to tens of seconds to accomplish and compared to fractions of a second to transition the processor from the third mode to the second mode or to the first mode.

The present invention provides and supports several different structures, mechanisms, and procedures for controlling the operational modes of the server modules and hence the processor or processors that may form or contribute to the operation of a server. Organizationally, the control may

reside in a separate Management Module, one or two of which Management Modules may be integrated into one of the inventive ISSU; or, may reside in one of the Server Modules which has been designated as a manager, supervisor, or master server module. Designation of a Server Module in this way involves providing the server module with the computer program software for receiving activity information from the server modules, for analyzing the activity information to determine from a power consumption perspective (or other predetermined perspective) which server modules should be operated in the several available modes (for example, Mode 1, Mode 2, Mode 3, and Mode 4 in some circumstances), and where the operation of more than one type of server is to be combined into a single server module (such as a multi-media server and a web page server) for gathering the content from the type types of servers onto the hard disk drive of a single server or group of servers. Note that when a particular server module is to serve as the master, that server may collect information on its own activity and be considered in the overall server and power management scheme. Of course, the server module acting as its own master will not typically place itself in either Mode 3 or Mode 4 as its continued operation is necessary to control other server modules. Where appropriate logic is provided to place the master in a power conserved state (such as Mode 3) and bring it out of that state, even the master may be placed into one of the Mode 3 operating states.

At a top level, the server modules each detect and report their activity to the manager (either the management module or the designated master server module). In some embodiments, the server modules are permitted to locally control their own operating mode, for example whether their own CPU (or CPUs if a multiple CPU server) is or should be operating in a Mode 1 or Mode 2. They will then also report not only their activity level but also the operating mode under which the reported activity was measured or detected.

At another level, the manner in which activity is detected is an issue. At yet still another level, the power management control policy or procedure, that is the control plan that regulates which server modules should be placed in which of the available modes to provide both the required (or desired) performance according to some measure and the required (or desired) power conservation. Those workers having ordinary skill in the art will appreciate, in light of the description provided here, that there are virtually limitless different policies for power management. Specific policies that optimize or near-optimize the combination of server performance and power conservation may be determined empirically during initial installation and operational phases as they will likely depend upon the content served, the variation of server loading as a function of time of day, advertising or promotions, average server loading, amount of over-provisioning, minimum quality of service requirements, power consumption of server modules versus content served, and other factors associated with server operation. The policies may also be modified according to the particular physical and/or electronic or logical structure of the servers. Even different CPU technologies may suggest different policies.

It may also be observed that such policies may be biased in favor of any one or combination of server operational factors. For example, operation and therefore the policy for control may favor power conservation even though there may be some impact on performance. Alternatively, the policy may favor absolutely maintaining a quality of service even if power conservation is somewhat sacrificed.

As general multi-power management policy it is observed based on analytical and empirical data, that there is a certain power consumption overhead associated with each server device and that it is therefore generally preferred to operate a minimum number of server modules at near their maximum output (Mode 1). When a single device approaches its capacity, other server devices are brought up from a Mode 3 to Mode 2 or Mode 1 operation. Frequently, the two servers then on line would each operate in Mode 2 until further performance is needed, at which time one would be brought to Mode 1 operation. This is merely an example scenario and many other alternative control strategies may be applied. Clearly, there is a bodies of knowledge for both open-loop and feed-back based control that may be used by those skilled in the art to optimize or near-optimize some weighted combination of performance and power conservation.

A server system configuration tool may be provided that allows a human operator to monitor system operation and power consumption and interact with the system and policy definition within the system to tune system performance. In the event that local government or regulatory agencies restrict power consumption or mandate power reduction, the policy may be altered to implement these requirements. In each of these situation, the system permits real-time dynamic uploading of the policies without taking an servers offline. In one embodiment, systems having two management modules are used effectively by off loading one management module to the other management module, updating the policies in the off loaded management module, and then placing the updated management module. In another embodiment, alternative policy schemes are preloaded in the management module (or designated master) so that it may switch automatically or under operator control as required.

In one embodiment of the invention, the computer system comprises a server for serving data or other content in response to a request. A hypothetical scenario in which a computer system, which may typically be but not necessarily be a portion of a larger network system having multiple server computers, transitions from a full power maximum performance operating mode to an off state in which the computer system neither performs operations no maintains state. The particular progression between states or modes may possibly but is unlikely to occur in a real computer system as it is more likely that certain modes will be skipped either to reduce power consumption when performance requirements are low or skipped when performance demand increases so as to elicit a higher performance operating mode than the next progression would provide. In general, the inventive system and method may provide for transitioning between an one state and any other different state. In some embodiments of the inventive system and method, not all of the modes described here will be present. Furthermore, other embodiments of the invention may provide for additional and different control. Furthermore, the description immediately below addresses control of the processor unit (e.g. processor or CPU) and logic circuits (frequently referred to as core logic or SouthBridge) associated with the processor unit. It should be understood that control of other components within the system, including for example hard disk drives, input/output ports, network interconnect circuits or cards, BIOS, video circuits, clock generators, voltage regulators, micro-controllers, memory, as well as other individualized logic circuit components may be independently or dependently controlled or controlled as groups. (See for example, Table III and the accompanying description for the manner in which some elements are controlled.)

It is initially assumed that the system is operating in Mode 1 having the highest processor unit (e.g. CPU) performance level and greatest power consumption of the available operating modes. The system is configured with operating system software (e.g. Microsoft Windows, Linux, Unix, Sun, or the like) and/or applications program software that include instructions for monitoring the occurrence or non-occurrence of an event.

It is noted that the Linux Operating system, such as the RedHat Linux operating system, may be more power conserving than other currently available operating systems. One reason for its power conservative features are the fewer number of instructions that need to be executed to accomplish tasks. Therefore while embodiments of the invention support all of the available operating systems, and may be adopted to support future operating systems, one embodiment utilizes the Linux operating system to achieve a higher degree of power conservation.

One such event that can be monitored and detected is the occurrence of execution of an idle thread. Another such event is the occurrence of some specified level of CPU processing capability availability that is derived from some enumeration or statistical evaluation of the idle thread or idle threads that are being or have been executed during some time period. Other events that may trigger a transition are described elsewhere in this specification. For purposes of continuing the description here, it is assumed that execution of idle threads is monitored and reported by a combination of an application program and the operating system, and that the number of idle threads being executed suggests that more performance is available than is needed and that power consumption may be reduced without sacrificing performance.

Control signals are then generated (either locally by the CPU or core logic, or globally by a separate power manager) that transition the system from Mode 1 to one of the Mode 2 operating modes. Mode 2 is generally characterized by having a CPU clock frequency that is less than the maximum rated CPU clock frequency, a CPU core voltage less than or equal to the rated maximum CPU core voltage, and core logic that operates at or substantially at the rated core logic clock frequency and core logic operating voltage. (This condition is also referred to as the Mode 2' operating mode.) By maximum rated CPU clock frequency is alternatively meant: (i) the clock frequency the CPU manufacturer has identified with this CPU model, (ii) the actual maximum frequency at which the CPU may be clocked, (iii) the maximum clock frequency that the CPU is operated within the system independent of what the CPU is capable of being operated at, (iv) or some similar or analogous measure. For example, if the CPU is marketed or sold as a 800 MHz Intel Pentium III, then the maximum rated CPU clock frequency is 800 MHz. If the maximum clock frequency at which the 800 MHz Intel Pentium III is operated in the system is 850 MHz, then the maximum rated frequency is 850 MHz.

It is also understood that there are gradations of performance (and power consumption) within the rubric of Mode 2 operation. A Mode 2'' operating mode is characterized by operation at both less than the maximum rated clock frequency and at less than the maximum rated core voltage. Mode 2 may be a single operating mode, or include a plurality of operating modes, having the general Mode 2 characteristic but providing for several different CPU clock frequencies and core voltage that at least support electrical device switching (transistor switching) or be selected to provide just adequate core voltage substantially matched to the clock frequency to provide reliable operation. For

example, at the Mode 2" operating mode, the CPU clock frequency and CPU core voltage are the minimum clock frequency and core voltage that are operable and supported by the CPU (where such minimum exists). Embodiments of the inventive system typically provide that core logic continue to operate at nominal rated levels where both the core logic clock frequency and core logic operating voltage are at or substantially at rated levels. In other embodiments, of the invention core logic circuit elements may also be power managed during Mode 2 operation by reducing clock frequency, operating voltage, or both.

The CPU clock frequency may be adjusted over a range of frequencies to match the amount of processing capacity to the tasks to be performed. Therefore, as the number of idle threads being executed in the CPU continue to increase indicating that productive tasks (such as retrieving data from a storage device, and sending such retrieved data to an I/O port or NIC for serving to a requester) are being performed within specified limits or some quality of service parameters, the clock frequency may be continually reduced.

At some time, however, the requirements placed on the system may become some low that at times there are no tasks to be performed. For example, on a computer network having a multiplicity of servers for serving stock market quotes and having sufficient capacity to handle worst case traffic in an active stock market, there is likely to be lots of over capacity of a national holiday where the stock markets are closed and there is little interest among investors. Under such conditions (actually likely under less strenuous conditions than these) the CPU within a computer system may complete all pending applications or user tasks and begin executing a system idle loop. Such an idle loop may initially or after some period of time cause execution of a CPU Halt instruction (or the equivalent) that causes the CPU clock to stop. This CPU halt instruction may be generated by the CPU itself or through some other internal or external agent or program. For example, a Microsoft Windows operating system or a Linux operating system are capable of generating an instruction to halt the CPU or processor. A halted or stopped CPU is one example of a Mode 3 operating mode, and more particularly a Mode 3' operating mode that is nominally characterized by a stopped or substantially stopped CPU clock, and a CPU core voltage that is less than or equal to the nominal maximum CPU core voltage and more usually at the minimum CPU core voltage that is necessary to maintain CPU register state and/or other CPU state. A CPU suspend state is another example of a different type of Mode 3 operation. Mode 3" may represent further power conservation by lowering the CPU core voltage to that just required to maintain state. This is treated as a separate sub mode because CPU core voltage need not be reduced as a result of the CPU halt command, and as stopping the CPU clock for a short period of time between execution of application tasks itself provides significant power savings without the design changes that may be required to also transition core voltage. Reduction of core voltage when the clock is stopped also generally has a smaller impact on power conservation than when the CPU is clocking. Some embodiments will also operate the CPU at the minimum clock frequency and minimum CPU core voltage as provided under a Mode 2 operation, and when executing the CPU halt instruction turn off the clock from that minimum value and maintain the core voltage at the voltage that supports the minimum clock. In this manner, the CPU may halted and resumed from halt by restarting the

clock and leaving the voltage alone. This scenario may be particularly effective when making rapid transitions between Mode 2 and Mode 3.

When it is determined that the CPU and computer system in which the processor is installed are not needed for some longer period of time, it is possible to provide additional power savings by reducing the power consumed by the core logic circuits or chips associated with the CPU. Where this additional level of power reduction is desired, the core logic clock frequency may be reduced to something less than the nominal or maximum frequency and in addition but optionally, the core logic voltage may be reduced so as to support that frequency. CPU and core logic state are maintained in each of the Mode 3 operating modes.

When the computer system is not needed for some longer period of time, the processor or CPU and at least a substantial portion of the core logic may be turned off. This is represented by Mode 4 operation which in one embodiment is characterized by having the CPU core voltage at zero, the CPU clock frequency at zero, most of the core logic circuits receiving no operating clocks or operating voltage. In some embodiments, the real-time clock may continue to operate and/or one or more circuits may remain active so that they may receive an external signal (such as a Wake-on-LAN) derived signal and be turned back on to resume operation in one of Modes 1, 2 or 3.

Note that in some embodiments, wherein if a portion or the entire system is operating in a reduced power consumption mode, such as one of the mode 3 operating modes, the manager or supervisor (such as a management module determines that server modules are dropping packets and that few or no idle threads are executing (indicating that the system has insufficient performance capability) then the supervisor or manager can send a packet to the fast Ethernet controller (or other circuit) associated with the server module to wake it up. This packet may be any packet identified to the ethernet controller (or other controller) to wake up the server. In one embodiment, the line or bus is monitored for special "ON" packet. In another embodiment, any packet received will turn it on. This wake up feature is advantageous as when the processor or CPU is not being clocked (such as in a Mode 3 operating mode) additional means are needed to wake it up to place it in a active mode that can process instructions and perform useful tasks, and in a server environment, the server according to embodiments of the invention will be connected to the ethernet and active. Other types of wake up or attention signals may alternatively be used.

When performance requirements increase, the computer system may generally transition from lower performance (and lower power consumption) modes to higher performance (and typically higher power consuming modes) according to rules, policies, algorithms, and/or control mechanisms provided in the system. Transitions may also occur dynamically. The events which trigger change or transition from one operating mode to another operating mode may typically be set and changed under programmatic software or firmware control. Various exemplary situations or events that trigger transitions are described elsewhere in this specification.

While a number of modes (Mode 1, Mode 2, Mode 3, and Mode 4) have been described in this example, it is noted that the inventive system, method, and computer programs do not require each of these modes or each of the submodes (e.g. Mode 3") within a mode. Furthermore, depending upon the configuration of the system, the set of rules or policies in place during operation, and/or the dynamics of operation

at the time an operating mode decision is to be made, for any single computer system, or group of computer systems, and their included processor, processing unit, or CPU, operation may transition between any two of the modes described. The examples provided here and the modes or states identified in the state diagrams are therefore illustrative rather than limiting.

By way of highlighting selected ones of the computer system (for example, server module computer system) operating modes, embodiments of several of these modes and submodes are now briefly described.

One embodiment of a first mode (Mode 1) comprises a mode in which the processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

One embodiment of a second mode (Mode 2) comprises a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

One embodiment of a second submode (Mode 2') further comprises a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

Another embodiment of the second submode (Mode 2'') further comprises a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

Another embodiment of a second submode (Mode 2''') further comprises a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage just sufficient to maintain switching circuits in the processor unit at the processing unit clock frequency, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

One embodiment of a third mode (Mode 3) comprises a mode in which the processing unit is operated at a slow but non-zero frequency processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage;

One embodiment of a third submode (Mode 3') further comprises a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage;

Another embodiment of a third submode (Mode 3'') further comprises a mode in which the processing unit is operated at a substantially zero frequency processing

unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

Another embodiment of the third submode (Mode 3''') further comprises a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than or equal to a maximum rated logic circuit operating voltage.

Another embodiment of a third submode (Mode 3''''') further comprises a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than a maximum rated logic circuit operating voltage.

Another embodiment of a third submode (Mode 3''''''') further comprises a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a substantially zero logic circuit clock frequency and at a logic circuit operating voltage that is just sufficient to maintain logic circuit operating state.

One embodiment of a fourth mode (Mode 4) comprises a mode in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage.

An embodiment of a fourth submode (Mode 4') further comprises a mode in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and the logic circuit is powered off by removing the logic circuit clock and by removing the logic circuit operating voltage or by setting the logic circuit operating voltage below a level that will maintain state, except that a real-time clock and circuit for waking the logic circuit and the processing unit are maintained in operation.

Another embodiment of a fourth submode (Mode 4'') further comprises a mode in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and the logic circuit is powered off by removing the logic circuit clock and by removing the logic circuit operating voltage or by setting the logic circuit operating voltage below a level that will maintain state, except that a circuit for waking the logic circuit and the processing unit are maintained in operation.

Some of the characteristics of these modes and submodes are listed in Table II. FIG. 4 provides an exemplary state engine state diagram graphically illustrating the relationships amongst the modes and identifying some of the transitions between states or modes for operation of an embodiment of the inventive system and method. Note that although the state engine may provide a path for directly or indirectly transitioning between any two modes or sub-



modes, in the interest of keeping the state diagram intelligible, the state diagram of FIG. 4 does not show all of the possible state or mode transitions possible.

Having described several power or energy consuming states or modes (or their opposite, power or energy conserving states or modes) as well as a situation in which a hypothetical computer system may transition between these modes, it will be appreciated that some procedure, mechanism, or policy is provided for the processor to self- or locally-control its own operating mode and hence its power consumption.

It is further noted that these operation modes may be utilized in different combinations and that any single system need not implement all of the operational modes. Therefore it will be appreciated that in the appurtenant claims, references to various modes, such as first mode, second mode, third mode, fourth mode, or the like, may refer to operating modes or states in a general manner as otherwise defined in the claims rather than to operating modes described in such terms in the specification. For example, in the claims where two operating modes are recited, such as first and second modes, such two modes may be any of the modes or states described, references, or suggested herein.

Heretofore, control of the operating mode of a plurality of processors or CPUs by a single supervisor or manager has not been known, particularly when the supervisor or manager is itself or includes a processor or CPU, and more particularly, it has not been known to provide this type of multi-processor power management in a multi-server system. This level of control is referred to herein as global control over a plurality of processors to distinguish from the afore described single processor or CPU power management.

It is noted that the inventive system and method also extend beyond any single-board computer systems having multiple processors configured therein. No such multi-CPU computers are known that provide power conservation features of the type described herein, and it is noted that in exemplary embodiments of the inventive system and method that each of the plurality of processors are located within separate PC-board mounted module. Embodiments of the inventive system and method are provided for which both local-control and global-control are provided. Such global control over a plurality of computers or appliances (each itself having either a single or multiple CPUs or processors) is not therefore limited to computers operating a servers.

Embodiments of the invention provide for detecting activity (or inactivity) in numerous ways, including but not limited to at least three different ways described herein. Detection may occur at the local level so that local control can be effected as well as optional detection at a global level. It is noted that in at least some embodiments, local detection of activity within each processor or CPU provides sufficient information to globally control the power consumption of a system having a plurality of processors or CPUs.

In one embodiment, an OSI model having a physical layer is used for activity or inactivity detection. In a second embodiment, a TCP/IP layer is used for this detection, and in a third embodiment the activity or inactivity detection occurs at the application layer. In a fourth embodiment, two or more of these activity detection and control techniques are combined.

TABLE II

Selected Example CPU and Core Logic Clock and Voltage Ranges for Various Exemplary Computer System Operating Modes.

CPU Mode	CPU Clock	CPU Core voltage	Core Logic Clock	Core Logic Voltage
1	≈max	≈max	≈max	≈max
2	<max	≤max	≈max	≈max
2'	<max	<max	≈max	≈max
2''	<max	<max	≈max	≈max
2'''	<max and >0	<max and sufficient to maintain switching rate (and CPU state)	≈max	≈max
2''''	≈min and >0	≈min and sufficient to support switching rate (and CPU state)	≈max	≈max
3	<max and ≥0 (typically = 0)	≤max and ≈min sufficient to maintain CPU state	≤max but >0	≤max but >0
3'	≈0	≤max and ≈min sufficient to maintain CPU state	≈max	≈max
3''	≈0	<max and ≈min sufficient to maintain CPU state	≈max	≈max
3'''	≈0	<max and ≈min sufficient to maintain CPU state	<max	≈max, or sufficient to support core logic clock freq.
3''''	≈0	<max and ≈min sufficient to maintain CPU state	≈0, except that generally RTC remains active	≈max
3'''''	≈0	<max and ≈min sufficient to maintain CPU state	≈0, except that generally RTC remains active	<max and ≈min sufficient to maintain logic state
4	=0	=0	most core logic circuits receive no operating clock	most core logic circuits receive no operating voltage
4'	=0	=0	core logic circuits receive no clock except for RTC and wake-up circuit	core logic circuits receive no voltage except for RTC and wake-up circuit

TABLE II-continued

Selected Example CPU and Core Logic Clock and Voltage Ranges for Various Exemplary Computer System Operating Modes.				
Mode	CPU Clock	CPU Core voltage	Core Logic Clock	Core Logic Voltage
4"	=0	=0	core logic circuits receive no clock except for RTC and wake-up circuit	core logic circuits receive no voltage except for RTC and wake-up circuit

One technique for detecting activity or inactivity in the physical layer uses idle thread detection. In certain operating systems prevalent in the late 1990's through 2001 provide a procedural "hook" through an operating system functional call or other programming construct that allows query of the operating system and generation of a response or report back to the requestor indicating how much idleness is present in the system, or more particularly how much idleness is present in the processor or CPU on which the operating system is executing. This operating system query may for example be made using an API function call that returns a value. Some exemplary alternative techniques and procedures for determining idleness in a system utilizes somewhat heuristic idleness detection algorithms, such an approach is described in U.S. Pat. Nos. 5,396,635, 5,892,959 and 6,079,025 by the inventor of the present application as well as in the other applications related thereto.

With reference to FIG. 5-12, several exemplary mode or state diagrams are illustrated. In these diagrams, a mode or state is represented by a circular node and a transition between two modes is represented by a directional line or arrow, the arrowhead indicating the direction of the mode transition. It is assumed for purpose of this discussion that the system may be in any one of three modes (Mode 1, Mode 2, or Mode 3) and a powered-off mode (Mode 4) (not shown).

Some systems, such as certain Transmeta Crusoe™ CPUs operate so as to provide a maximum CPU core voltage and maximum CPU clock frequency in a Mode 1 type operation and a plurality of levels (15 levels) in a Mode 2 type operation, the Transmeta Mode 2 operation consuming less power in fifteen of its operating levels than in the sixteenth operating level. Each of these fifteen lower power consuming levels at which the CPU core voltage and CPU clock frequency are less than their nominal rated maximum are considered to be Mode 2 operating states as the processor operates in at most one of the states at any given time and each separately qualifies as a Mode 2 operation relative to the maximum performance Mode 1 state and CPU suspend Mode 3 state. A mode state diagram for the Transmeta Crusoe LongRun™ CPU operation is illustrated in FIG. 7.

It is also noted that the Intel SpeedStep™ technology involves the same or similar three modes of operation. The Intel SpeedStep provides for a fully on mode running at maximum clock frequency and maximum CPU core voltage, it also has a reduced state in which frequency and voltage are reduced relative to maximum, and a suspend state. During normal operation such as for an AC-line powered notebook computer, the CPU clock frequency and CPU core voltage are at their rated maximum values. However, in at least one notebook computer made by IBM (IBM ThinkPad T21) a user may enable an optional power saving policy for battery powered operation and for AC-line powered operation in which the CPU clock frequency and the CPU core voltage

are reduced to save power and lengthen battery life. These power saving policies also control hard disk drive, display brightness, and the operating condition of other internal circuits and peripherals.

Each of FIG. 5-12 shows a first mode (Mode 1), a second mode (Mode 2), and a third mode (Mode 3). A fourth mode (Mode 4) represents a processor or CPU that is powered down or in an Off state and is not shown. Various mode transitions are supported by the inventive system and method. Conventionally, the transitions between and among the three modes were controlled locally (though such terminology was not used for such conventional systems because there was no global control to contrast with) because all or substantially all control was provided either within the CPU or by chips, logic, or other circuits associated with the single computer or PC-board on or in which the CPU was located. In aspects of the present invention, global control is exercised over the operating modes of a plurality of the processors or CPUs, and some degree of local control is or may optionally be provided. The manner in which the transitions are controlled locally and globally are described in greater detail elsewhere in this specification.

Recall that in single processor or single CPU systems, Mode 1 and Mode 2 represent active work producing operating states, a non-zero frequency processor clock signal causing the switching of transistor or other circuits that permit instruction execution. Therefore, in single processor systems, particularly in notebook computer systems operating from finite energy sources (e.g. battery), the systems occupy most of the time they are "on" in a Mode 1 condition (or Mode 1-like condition) or in a Mode 2 (or Mode 2-like condition). Operation in a Mode 3 condition does not provide any productive work so that if the user were to perform any reasonable amount of work using the device containing the power managed processor or CPU, there is little power savings that would be achieved during useful work.

In FIG. 5-12 the following notation is adopted. Each transition indicating arrow is labeled with either an "L" to indicate local control, a "G" to indicate global control, or an "LG" meaning that the transition may be controlled by either or both local control or global control. In addition, transitions from Mode 1 to Mode 2 are labeled "A" and transitions from Mode 2 to mode 1 are labeled "A'". In analogous manner, other transitions are labeled as B, B', C, and C'. This notation will be useful in describing the differences between conventional systems and method and the present invention.

With respect to FIG. 5, there are shown locally controlled transitions between Mode 1 and Mode 2 (A and A') and between Mode 2 and Mode 3 (B and B'). For recent power management schemes, the A and A' transitions would normally be expected to occur with reasonable frequency during use of the notebook computer, and the B and B' transitions with lower frequency, under the assumption that

the user will typically either be using the computer (A and A' transitions) or power it off (Mode 4), so that B and B' transitions will be less frequent. It may also be expected that the B' transition may be less frequent than the B transition, as computer makers may typically transition directly to Mode 1 from a Mode 3 (C' transition) when there is suddenly a need to wake up the CPU from a suspend type state. It is noted that for embodiments of the present invention, the B and B' transitions may be frequent to very frequent, particularly when the 3rd mode is the Mode 3' state in which only the CPU clock is halted and all or most other system clocks remain operational. The Mode 3' to Mode 2 (or Mode 1) and the Mode 2 (or Mode 1) to Mode 3' transition can occur very rapidly and because of the high CPU clock frequency and the number of switching circuits present in modern CPUs can yield considerable power or energy savings. Embodiments of the invention may also provide that a system operating in Mode 3' (CPU clock stopped or slowed significantly) may also further transition to a Mode 3'' (CPU and other clocks stopped or slowed significantly) under specified conditions.

FIG. 6, illustrates an operating scenario under which the processor or CPU is maintained in an active state and only the A $\rightleftharpoons$ A' transitions occur under local control. The B $\rightleftharpoons$ B' and C $\rightleftharpoons$ C' transitions are illustrated in dashed lines.

FIG. 7, illustrates a similar operational scenario wherein the processor or CPU may transition to any one or sequentially through a plurality of Mode 2 states. This operational scenario is similar or the same as the scenario under which the Transmeta Crusoe processor may operate.

The inventive architecture, system, device, and method may be operated in a fundamentally different manner, using either only global control or using a combination of local and global control, to alter the operating mode of a plurality of processors or CPUs. Variations on this power management scheme are now described relative to FIG. 8-12.

In FIG. 8, the Mode 1 to Mode 2 A $\rightleftharpoons$ A' transitions are locally controlled. For example, in the Intel SpeedStep™ CPUs the A $\rightleftharpoons$ A' transitions are controlled using control mechanisms provided by Intel on their CPU chips that permit a system designer to issue a command to the CPU to transition it from Mode 1 to Mode 2 under an identified condition and from Mode 2 to Mode 1 under a second identified condition. Similarly, the Transmeta Crusoe CPUs implementing their LongRun technology would transition from Mode 1 to a selected one of a plurality of Mode 2 states, and from that Mode 2 state (or a different Mode 2 state) to Mode 1, under identified conditions. These conditions are known in the art, available from Intel or Transmeta, or from Intel, AMD, or Transmeta computer manufacturer OEMs, and not described here in greater detail.

While the conventional systems and methods may permit the B $\rightleftharpoons$ B' transitions and/or the C $\rightleftharpoons$ C' transitions under local or self-control within a processor or CPU (or within circuitry associated with a CPU on a common mother board or other platform or enclosure), embodiments of the inventive system and method preclude such local or self-control. Rather, a manager or supervisor (see description of manager or supervisor capabilities and implementations elsewhere in this specification) only may globally manage the B $\rightleftharpoons$ B' transitions and/or the C $\rightleftharpoons$ C' transitions under a global control scheme. Global control in this manner is illustrated for example, in the state diagram of FIG. 9.

In yet another embodiment of the invention, depicted in the FIG. 9 state diagram, Mode 2 operation is not supported and there are no A $\rightleftharpoons$ A' transitions or B $\rightleftharpoons$ B' transitions. It is observed that operating only in Mode 1 or Mode 3 would not

represent a generally useful power management scheme for a single processor or CPU system because Mode 1 operation is a full power active mode and Mode 3 is power conserving but inactive mode. Therefore, there is little power savings that would result where CPU or processor loading is sufficient to keep the processor or CPU out of Mode 3. Significantly, systems or power management policies providing only C $\rightleftharpoons$ C' transitions for single CPU systems (or for any processor or CPU systems) do not seem to exist in the computer industry.

On the other hand, this operating scheme is viable and presents significant power conservation features for multi-processor or multi-CPU architectures, particularly in the server environment where some or significant over-provisioning of server capacity is the norm and where the server suite may typically operate at from twenty to fifty percent of maximum capacity. As described in greater detail elsewhere in this specification, in the inventive Integrated Server System Unit (ISSU) a plurality of server modules, each having a processor, are integrated into a single enclosure and coupled for communication by various in-band and out-of-band bus and interconnection links. A manager or supervisor is provided (for example, in the form of a Management Module or designated Server Module operating as the manager or supervisor) that collects and/or analyzes CPU "activity" (where activity is defined broadly as described elsewhere in this specification) and generates control signals that maintain or alter the operating mode of individual Server Modules or identified groups of such Server Modules. While the primary control is over the processor or CPU within these Server Modules, it is noted that other circuits or components, such as for example, display, hard disk drive, and other circuits and/or peripherals may be similarly controlled by the same or different control signals.

Servers, server systems, or so called server farms generally designed and implemented with significant capacity over-provisioning. Reasons and rationale for such over-provisioning is known in the art and therefore described only briefly here. Providing a positive first visit Internet web experience and maintaining a quality of service (QoS) is important for developing and maintaining clients, customers, or other visitors to a web site. Content must be served within a reasonable period of time, on a first visit and on subsequent visit, or visitors will not return. While the quality of service may be permitted to vary somewhat by time of day and/or season, the reasonableness standard still applies, and normally it is best to maintain a very high quality of service all the time. Paramount in this goal would be to serve content such as web pages, streaming video, or cached content, without delay. Even during time periods (time of day, season, event driven) where web traffic and the amount of content that need be served by a server is likely to increase, sufficient server capacity must be in place. Over provisioning by at least 30% or so is typical, and frequently 100%-500% or more over-provision or over-capacity may be provided.

This moderate to significant over-provisioning is accepted by the server community as a necessary cost item, both in terms of the cost to purchase and maintain the equipment, the cost to power the equipment, the cost to cool or remove the heat generated by the equipment, and the negative impact on equipment longevity as a result of continuous operation.

Conventional server systems have not been power managed as there has been a philosophy that if the equipment is there it should be operated at maximum speed so as to serve content or respond to other requests as rapidly as possible.

Conventional server units within a rack of server units have been to the inventor's best knowledge maintained in an always on always ready to serve mode. More recently, there has begun to be some appreciation that power saving features provided in commercial personal computers might result in some power conservation benefits. At most these recent ideas have concentrated on the Mode 1 to/from Mode 2 ( $A \rightleftharpoons A'$  transitions) based on the Intel SpeedStep™, Transmeta Crusoe LongRun™, or other similar technologies. This local self-control by each processor provides some energy conservation but does not provide the conservation of the inventive system and method.

One of the Transmeta Crusoe Model chips operates at 533 MHz and 1.6 volts when in Mode 1 and at 300 MHz and 1.2 volts when at its slowest CPU clock frequency and lowest CPU core voltage in Mode 2. (Note that these operating parameters are nominal and subject to change by their manufacturer from time to time as products change, even within a particular product model or family.) Recall that to a general approximation  $P \propto K_1 C f v^2 + K_2$ , where  $P$ =power consumption,  $f$  is clock frequency,  $v$ =CPU core voltage,  $C$ =capacitance,  $K_1$  is some multiplicative proportionality constant, and  $K_2$  is some additive constant that represents the small power consumed by a circuit when operating voltage (e.g.  $V_{cc}$ ) is applied but the CPU or processor clock is turned off (e.g. 0 MHz clock, or very slow clock). While these values may change for different CPU designs and chip sets it will be clear that the savings in transitioning from a 1.6 volt/533 MHz operation to a 1.2 volt/300 MHz operation is modest as compared to transitioning from a 1.6 volt/533 MHz operation to a 1.2 volt/0 MHz operation. Operation with a CPU core voltage that is equal to that of the CPU clock slowed Mode 2 or an even a lower CPU core voltage than that needed to maintain a 300 MHz clock switching may be used during Mode 3 operation when only CPU register and memory contents or status need be maintained.

It will therefore readily be appreciated in light of this description that operating a multi-server system where at least global control of the operating modes of a plurality of CPUs (and optionally other circuit elements of the servers) will yield significant power conservation benefits. Furthermore, in some operational situations combining Mode 1 to/from Mode 2 ( $A \rightleftharpoons A'$  transitions) either locally controlled or globally controlled may add even further power conservation features.

FIG. 8 illustrates the state transition for an inventive embodiment in which  $A \rightleftharpoons A'$  transitions are controlled locally, and  $B \rightleftharpoons B'$  and  $C \rightleftharpoons C'$  transitions are under the control of a global manager. FIG. 9 illustrates the state transition for an alternative inventive embodiment in which the processor or CPU only operates in either Mode 1 or Mode 3 and not in Mode 2 so that  $A \rightleftharpoons A'$  and  $B \rightleftharpoons B'$  transitions are prevented from occurring (such as by, disabling a feature provided with a chip, de-configuring power conservation features, or providing the manager with the ability to otherwise prevent such transitions), and  $C \rightleftharpoons C'$  transitions are under the control of the global manager.

FIG. 10 illustrates the state transition for yet another alternative inventive embodiment in which the processor or CPU only operates in any of Mode 1, Mode 2, or Mode 3 and while the  $A \rightleftharpoons A'$  transitions occur under local control, the  $B \rightleftharpoons B'$  transitions are prevented from occurring, and  $C \rightleftharpoons C'$  transitions are under the control of the global manager. In this embodiment, therefore, the transition to Mode 3 therefore only occurs directly from Mode 1 and never from Mode 2. In yet a further embodiment, illustrated in FIG. 11, the  $A \rightleftharpoons A'$  transitions occur under local control and the  $B \rightleftharpoons B'$

transitions occur under global control, and where  $C \rightleftharpoons C'$  transitions do not occur. FIG. 12 illustrates the mode transitions in a further embodiment, where each of the  $A \rightleftharpoons A'$ ,  $B \rightleftharpoons B'$ , and  $C \rightleftharpoons C'$  transitions may occur according to predetermined power management policies and where each separate possible transition may be under either local and/or global control according to the predetermined policy or power management procedure or algorithm. The policy, procedure, or algorithm may also disable certain states of transitions statically or dynamically, and may cause certain of the server modules or other CPU or processor based devices into a powered off (Mode 4) and back to any of the powered on modes.

FIG. 13 illustrates that for a system having a plurality of processor or CPU based devices, the CPU or processor within any particular device (such as server modules) may be in different states at different times under the direction of an device-local control, a system supervisory global control, or a combination of the two. The shaded mode circles indicate the current mode and the mode transitions, though not shown, may be any of those already described relative to the other inventive embodiments.

In light of the above description, it will be appreciated that the inventive system and method extends earlier power management structures, architectures, and methods by the same inventor Henry T. Fung (such as are described in U.S. Pat. Nos. 6,115,823; 6,079,025; 5,987,614; 5,961,617; 5,892,959; 5,799,198; 5,758,175; 5,710,929; and 5,396,635, herein incorporated by reference) to multi-server or multi-node architectures.

These existing power management patents include innovative systems, architectures, and methods for saving or conserving energy or power within a single system by using one or more of several power management schemes, including, but not limited to the following schemes: (1) Detection of the idle activities by monitoring I/O activities or execution of a predefined code thread. (2) Reduction of power consumption by lowering (or stopping) various clock frequencies or removal of power (operating voltage) to different components within the system. (3) While in a power saving mode, continuing to monitor the occurrence or non-occurrence of a second predefined event or activity and entering a deeper power saving mode in response to the second predefined event or activity detection. Note that although certain events, activities, and/or indicators are referred to predetermined, such events, activities, or indicators may be dynamically determined during operation as well as determined in advance.

The present Multi-Server Power Management scheme extends these earlier techniques, augments them, and introduces entirely new features and capabilities. Five particular innovations are set forth below, however, it will be apparent that the invention described herein is not limited only to this set of features and capabilities.

First, power management of the network devices including the server modules can occur at different OSI levels and be extended beyond the physical layer. In particular, the detection of server activity whether measured by idle activities or other means may occur at the physical layer but is advantageously extended beyond the physical layer to the network layer (for example, to the TCP/IP layer) and to the application layer. For example, at the physical layer, the number of CPU idle threads within a fixed time period may be detected or measured, or, some type of statistical evaluation of CPU idleness may be determined. As one numerical example, if the CPU is idle 80% of the time while in a particular operating mode such as Mode 1, it is clear that this

much processing performance is not required and the CPU performance may therefore be adjusted downward to save power. If we assume in a simple case that a Mode 2 operation reduces the CPU clock speed by a factor of  $\frac{1}{4}$  over the Mode 1 clock speed, then the CPU will only be able to process  $\frac{1}{4}$  of the instructions in the same period of time, however, this is sufficient given the 20% loading (80% idleness) the CPU is experiencing. Therefore, based on this idleness detection, significant power savings are realized. Alternatively or in addition, if for example, under the same scenario there is a group of ten network server devices that are being managed as a single logical group or image, eight of them may be put into an inactive but powered on Mode 3, and the other two network server devices placed in a Mode 1 operating state running at a 100% performance level.

Power management may also or alternatively occur based on detection at the TCP/IP layer (or equivalent layer where a protocol other than TCP/IP is implemented). Under this detection and control model, CPU performance is monitored relative to the handling of TCP/IP packets. CPU performance level is lowered, such as by reducing CPU clock frequency (desirably accompanied by a reduction of CPU core voltage) until packets start dropping, and then increasing performance so that packets are not dropped and to provide an operating margin. The initial reduction and subsequent increase in CPU or server performance may be accomplished by altering the operating mode of individual servers or by adjusting the aggregate characteristics of a group of servers to provide the aggregate performance required. It is noted that where communications channel bandwidth limits the performance of a server, there may be advantage to reducing the performance level of the server to just satisfy the bandwidth limitation and thereby conserve power in the server.

At the application layer, the activity monitoring or detection may for example involve measuring the number of times a specific port address is or has been requested within a fixed time period. This determination or measurement may be accomplished, for example, by using a SNMP agent. In response to this measurement, an appropriate number of servers each operating at an appropriate performance level (Mode 1 or Mode 2) are provided to meet the performance requirement for each application. The rest of the servers are placed in a highly power saving state (Mode 3 such as Mode 3' [e.g. CPU clock halted] or Mode 3'' [e.g. CPU and other logic clock stopped], or Mode 4). The policies for selecting the number of active servers and their operating mode are described elsewhere in this specification. Recall that different application types may use different rules or policies to determine the server CPU performance and power conservation requirements.

Second, power management is extended beyond a single processor or CPU and in particular is extended beyond a single server (independent of the number of processors it may contain) to multiple servers across an entire network. It will be appreciated that this multi-server power management capability may be provided either with discrete servers or with the particular embodiment of the Integrated Server System Unit (ISSU) or Integrated System Server architecture generally.

Third, activity information created by any one server (or server module in the ISS scheme) is accessible to a designated supervisor via standard networking protocol. This supervisor is frequently referred to as the master, the capabilities of the master residing for example in an ISS Management Module or an ISS Server Module, though the

particular location or processor responsible for accessing and utilizing the activity information for the servers is not critical to the power management. In preferred embodiments of the invention, the supervisor or master capabilities reside in one or more management modules, and in an alternative embodiment, the supervisor or master capabilities reside in a designated or selected one of the server modules.

Fourth, servers can be reconfigured to run a specific application (e.g. web, streaming media and email) based on a certain load distribution requirement or requirements existent at the time upon receiving commands from a designated supervisor or master. Advantageously, this feature will provide or support operation at three or more power consumption levels, including a first full power mode (full CPU core voltage and normal maximum CPU clock frequency), a second mode consuming less power than the first mode in which either the CPU core voltage or the CPU clock frequency or both are reduced from the first mode, and a third mode in which the CPU is substantially inactive and consumes less power or energy than the second mode. In one embodiment, this third mode provides a CPU core voltage to maintain state and either stops the clock or maintains the clock at a very low frequency (for example, 1 Hz to a few hundred Hz) so that the CPU is effectively inactive.

Fifth, allowing any number (including none, one, many, or all) of servers across the entire network to go in and out of a 3rd power consumption mode directly from a first mode (Mode 1) without going through another intermediate power saving mode upon receiving commands from a designated master. This third power consumption mode (Mode 3) may for example include a mode where the processor or CPU is powered at some level but substantially inactive from the standpoint of executing commands or serving content, and memory associated with the CPU is refreshed. This third mode may be further broken down into a mode in which only the CPU clock is stopped (Mode 3') such as may occur when a Halt instruction is executed, and into a deeper power savings mode in which the CPU clock is stopped and other clocks are also stopped (Mode 3''). It is noted that in a typical implementation, the real-time clock (RTC) will generally run all the time so that certain system timing events and alarms can be maintained. The third power saving mode may also or alternatively be a powered down mode (Mode 4), however, such operation is somewhat undesirable unless it is anticipated that the powered down (Mode 4) server module will not be needed for some appreciable period of time as a delay is associated with bringing the CPU and the server module within which the CPU is located back on line. The Mode 4 operation may therefore only be used when the Mode 4 operation is expected to continue for several seconds, minutes, hours, or longer periods of time. It will be appreciated that in the third power saving mode, the CPU clock (and or other clocks in the system) may be either off entirely or running at a very low rate (such as for example 1 Hz, 10 Hz, 100 Hz, 1 KHz, or some other value that is small in comparison to the nominal frequency (for example, typically in the 100 MHz to 2 GHz range) of the processors used for such servers. It will be appreciated in light of the description provided here, that the invention provides for direct transition between a full or substantially full power mode and an inactive or substantially inactive mode. Although, this power mode transition would be much less useful for battery-powered portable applications for notebook computers or PDAs because of the desirability of maintaining some activity such as when typing into a word processor, this transition scheme extremely useful in a multi-server environment, where each of a plurality of

servers can serve the same content and it is desired to reduce the number of active servers while maintaining sufficient ability to satisfy quality of service requirements or otherwise maintain operation with a subset of the total set of servers.

These five innovations (as well as others) may of course be combined in various ways to provide even greater synergism. For example, the first described innovation extending the detection of idle activities beyond the physical layer to the network layer and/or to the application layer, may readily be combined with the fourth described innovation wherein the servers can be reconfigured to run a specific application based on a certain load distribution requirement or requirements existent at the time upon receiving commands from a designated supervisor or master.

This combination may also be extended according to the second described innovation to include multiple servers across an entire network, independent of whether the servers are discrete or integrated ISSU-based server modules. This latter combination may be further enhanced by also implementing the third described innovation to provide that activity information created by any one server (or server module in the ISS scheme) is accessible to a designated supervisor or master via standard networking protocol.

In yet another embodiment, the fifth described innovation that provides for any number of servers is a system having a plurality of servers to transition directly from a full performance 1st mode to an inactive 3rd mode. This scheme generally representing a non-useful power management scheme when applied to any single computer or server, but providing considerable benefit when the plurality of servers are managed in combination to provide a desired level of performance and power consumption savings.

Table III describes the behaviors of selected component inside an exemplary computer system, such as a computer system configured as a server module, at the different power management modes (Modes 1, 2, 3, and 4) according to one embodiment of the invention. This embodiment implements somewhat different power management policies than the embodiment described relative to Table II and also addresses the manner in which certain other peripheral devices or other components may be power managed. The mode descriptions are therefore generically similar but the detail or submode descriptions differ somewhat, but such differences are semantic and each of the modes and submodes described in any of the embodiments are within the scope of the inventive system, apparatus, computer program, and method.

In this embodiment's first mode (Mode 1) the processor or CPU functionally able to execute instructions for operating system and application programs; CPU activities are monitored, and the internal CPU clock frequency and CPU core voltage may be lowered if activity level of the CPU falls below some threshold (predefined or dynamically determined threshold). The voltage regulator is set to deliver the maximum (or specified nominal) CPU core voltage, the clock generator, RAM, hard disk drive, core logic, NIC, BIOS, and Real-Time Clock (RTC) are ON. The video may independently be controlled to be on or off and may even be absent from the system as video signals frequently are not needed for server systems, except in some cases of set-up or service. A microcontroller ( $\mu$ C) is operative and remains in continuous communications with the Management Module (or with an different Server Module designated or selected to operate as a manager or supervisor.

TABLE III

Exemplary behaviors of selected components inside a computer system (e.g. server module) at the different power management modes according to one particular embodiment of the invention. Other embodiments support alternative or additional modes and transitions between modes as described for example in Table II.			
	1 <sup>st</sup> Mode	2 <sup>nd</sup> Mode	3 <sup>rd</sup> Mode
CPU	1) Execute applications 2) Monitor CPU activities 3) Lower internal CPU clock frequency and voltage if activity level falls below a pre-defined threshold (go to 2 <sup>nd</sup> mode)	1) Execute applications 2) Monitor CPU activities 3) Go to the 1 <sup>st</sup> mode if activity level rises above a pre-defined threshold 4) Go to the 3 <sup>rd</sup> mode after receiving commands from an external master via standard network protocol (In Band communication)	1) CPU is in very low power state 2) Return to 2 <sup>nd</sup> mode or 3) Return to 1 <sup>st</sup> mode
Voltage Regulator	CPU core voltage is set to maximum	CPU core voltage is set to less than maximum	CPU core voltage is set to be equal to or less than core voltage in 2nd mode setting.
Clock Generator	ON	ON	Stop most (or all) clocks. For example, may stop only CPU clock, or may stop CPU and other clocks. (Usually RTC is not stopped.)
RAM	ON	ON	Suspended (refresh only)
Hard Disk	ON	ON	Suspended after receiving commands from the CPU
Core Logic	ON	ON	Suspended after receiving commands from the CPU or signal from $\mu$ C
NIC	ON	ON	Suspended after receiving commands from the CPU or turning off NIC Clk. Send resume signal to core logic after a

TABLE III-continued

	1 <sup>st</sup> Mode	2 <sup>nd</sup> Mode	3 <sup>rd</sup> Mode
Exemplary behaviors of selected components inside a computer system (e.g. server module) at the different power management modes according to one particular embodiment of the invention. Other embodiments support alternative or additional modes and transitions between modes as described for example in Table II.			
Video	ON/OFF	ON/OFF	predefined packet is received (e.g. Wake-On-LAN) Suspended after receiving commands from the CPU or turning off Video Clk
BIOS	ON	ON	Suspended
RTC	ON	ON	Send resume signal to the core logic after alarm expire
micro-controller (UC)	Continuous communications with the management module.	Activate the suspend signal of the core logic PMU (Out of Band) after receiving commands from the management module and causes the entire system to enter the 3 <sup>rd</sup> mode	Send resume signal to core logic after receiving commands from the management module

In Mode 2, the CPU still executes operating system and application program instructions, CPU activity is still monitored, and if the activity level rises above some predetermined or dynamically determined threshold (or according to some other rule or policy) the CPU enters Mode 1 operation, but the CPU enters Mode 3 in response to receipt of Mode 3 entry commands received from a manager of supervisor. These Mode 3 entry commands may generally be received from an external master via standard in-band network protocols. Recall that in Mode 2 the voltage regulator that supplies CPU core voltage is set to less than maximum core voltage. As in Mode 1, the clock generator is on but will (in preferred embodiments) deliver a lower frequency clock signal, and RAM, hard disk drive, core logic, NIC, BIOS, and Real-Time Clock (RTC) are ON. The Video may independently be controlled as in Mode 1. A microcontroller ( $\mu$ C) is operative in Mode 2 to activate a suspend signal of the core logic power management unit or PMU (Out of Band) after receiving commands from the management module (or server module acting as a designated master or manager) and causes the particular server and/or multiple servers within the system to enter the 3rd mode.

Is In Mode 3, the CPU is placed in a very low power consumption state and can return to Mode 1 or to Mode 2 upon the occurrence of some predetermined condition such as are described elsewhere in this specification. The voltage regulator that provides CPU core voltage is set to a voltage equal to or less than the core voltage in Mode 2 to thereby save power over that consumed in either of Modes 1 or 2. The clock generator is also stopped so that power consumed switching devices is substantially eliminated. (It is noted that in an alternative embodiment, the clocks in Mode 3 may be operated as a very slow rate, for example a few Hz to a few hundred Hz, or some other low clock frequency relative to the normal clock rate of the CPU.) RAM is suspended (that is the memory contents are refreshed only), the Hard Disk drive or drives are suspended after receiving commands from the CPU (or other commands to spin down and go into a suspend state). The core logic is also placed into a low power consuming suspend state after receiving a command from the CPU or signal from the micro-controller.

25

Mode 3 operation also provides for suspension of the network interconnect card or circuit (NIC) after receiving commands from the CPU or turning off the NIC clock. (Note that a resume signal may be generated and sent to the core logic if a predefined packet is received, such as for example, a Wake-On-LAN signal.) The BIOS is suspended, and the RTC may send a resume signal to the core logic after a RTC alarm expires. The microcontroller continues to monitor communications with the management module or other designated master so that it may send a resume signal to the core logic after receiving commands directing this type of action from the management module or other designated management master. If the Video was on prior to entering Mode 3, the Video is suspended after receiving commands from the CPU or the Video Clock signal is stopped or turned off, and if it was off it remains off.

While much of the description herein has focused attention on performance and power management of the processor, CPU, core logic, and other logic circuits within a computing device or system, or other information instrument or appliance having such processor and/or logic, it should be understood that the dynamic power management and dynamic workload management is not only limited to such systems or components. More particularly, the inventive dynamic power management system, method, architecture, procedures, and computer programs may also be applied to a diverse set of electrical and electronic components including components commonly referred to as computer peripherals. Application of the principles described herein therefore have the potential of reducing power consumption and prolonging component life to such devices and systems as video monitors, hard disk drives or other storage systems or devices, printers, scanners, cameras, other network devices and circuits, industrial tools and systems, and a myriad of other systems and devices.

Aspects of the invention though often described in the context of processors, CPUs, network devices, servers, and the like; have particular benefits relative to power and every conservation when applied to server farms where large quantities of energy are conserved directly as a result of lower power operation without performance sacrifice as

well as energy conserved as a result of higher density and lower facilities space and cooling requirements.

#### Additional Embodiments

Having described numerous embodiments of the invention, it will be apparent to those workers having ordinary skill in the applicable arts that the invention provides a great variety of innovations. Attention is now directed to highlights of the manner in which selected aspects of the invention and innovations may be used either separately or in combination to provide particularly desirable and advantageous utility. Although these highlighted groups of innovations and particular embodiments with each group are particularly useful, the inventions and innovations described in this specification and the drawings are not limited only to the embodiments highlighted or otherwise described or identified below. Within each group of innovations, the selected embodiments are, for convenience of notation, referred to by embodiment numbers surrounded by parentheses. These numbers refer to embodiments within a particular group of innovations and are reused for the different groups of innovations.

In a first group of innovations, the invention provides various embodiments associated with System, Method, and Architecture for Dynamic Server Power Management and Dynamic Workload Management for Multi-server Environment.

(1) A computer system comprising: a plurality of server computers each having at least one processor and an activity monitor identifying a level of activity indicator for the at least one processor; each of the server computers being operable in: (i) a first mode having a first maximum performance level and a first power consumption rate, (ii) a second mode having a second maximum performance level lower than the first maximum performance level and a second power consumption rate lower than the first power consumption rate, and (iii) a third mode having a third maximum performance level lower than the second maximum performance level and a third power consumption rate lower than the second power consumption rate; and a power manager: (i) coupled to each of the server computers and receiving the level of activity information from each of the plurality of computers; (ii) analyzing the plurality of received level of activity information; (iii) determining an operating mode for each of the server computers selected from the first mode, second mode, and third mode based on the analyzed activity information and predetermined policies; and (iv) generating commands to each of the plurality of server computers directing each of the plurality of server computers to operate in the determined operating mode.

(2) A computer system comprising: a plurality of computers each having at least one processor and an activity monitor identifying a level of activity indicator for the at least one processor; each of the computers being operable in: (i) a first mode having a first maximum performance level and a first power consumption rate, and (ii) a third mode having a third maximum performance level lower than the first maximum performance level and a third power consumption rate lower than the first power consumption rate; and a power manager: (i) coupled to each of the computers and receiving the level of activity information from each of the plurality of computers; (ii) analyzing the plurality of received level of activity information; (iii) determining an operating mode for each of the computers selected from the first mode and third mode based on the analyzed activity information and predetermined policies; and (iv) generating

commands to each of the plurality of computers directing each of the plurality of computers to operate in the determined operating mode.

(3) The computer system in embodiment (2), wherein: each of the computers further being operable in (iii) a second mode having a second maximum performance level intermediate between the first maximum performance level and the third maximum performance level and a second power consumption rate intermediate between the first power consumption rate and the third power consumption rate; and the power manager further determining an operating mode for each of the computers selected from the first mode, the second mode, and the third mode based on the analyzed activity information and the predetermined policies. (4) The computer system in any of embodiments (2 or 3), wherein: the computers comprise servers. (5) The computer system in any of embodiments (2, 3, or 4), further comprising a power manager computer providing the power manager. (6) The computer system in any of embodiments (2, 3, or 4) wherein a selected one of the plurality of computers designated as a master providing the power manager. (7) The computer system in any of embodiments (2 or 3), wherein the activity monitor comprises an activity monitor that monitors an activity selected from the set of activities consisting of: a program application layer activity, a network layer activity, a physical layer activity, and combinations thereof. (8) A system as in embodiment (7), wherein at the physical level the number of processor idle threads executed within a predetermined period of time are measured to determine processor loading and the processor performance is adjusted to by altering the operating mode to substantially match the level of processor loading. (9) The computer system in embodiment (2), wherein the activity monitor comprises a network layer activity monitoring TCP/IP protocol data packets; and processor performance is incrementally lowered by the power manager using the mode control until data packets start dropping indicating that the processor performance is at the limit of adequacy and then increasing the processor performance by a specified increment to act as a safety margin to provide reliable communication of the packets. (10) The computer system in embodiment (7), wherein the application layer activity monitor comprises monitoring use of a port address within the computers, the monitoring including counting or measuring a number of times a specific port address is being requested within a predetermined period of time, and in response to that counting or measurement, placing a sufficient amount of computer performance to meet the performance requirement for each application requesting the port address. (11) The computer system in embodiment (7), wherein the application layer activity monitor comprises monitoring use of a port address within the computers. (12) The computer system in embodiment (7), wherein the network layer activity monitor comprises monitoring use of a TCP/IP protocol within the computers. (13) The computer system in embodiment (7), wherein the physical layer activity monitor comprises monitoring the execution of idle threads within the computers. (14) The computer system in embodiment (7), wherein the physical layer activity monitor comprises monitoring counting activities having particular activity values within the computers. (15) The computer system in embodiment (3), wherein: the first mode operation is characterized by operating the processor at a first processor clock frequency and a first processor core voltage, the second mode operation is characterized by operating the processor at a second processor clock frequency and a second processor core voltage, and the third mode operation is characterized by operating



the processor at a third processor clock frequency and a third processor core voltage; the second mode of operation being further characterized in that the second processor clock frequency and the second processor core voltage in combination consuming less power than the first processor clock frequency and the first processor core voltage in combination, and the third processor clock frequency and the third processor core voltage in combination consuming less power than the second processor clock frequency and the second processor core voltage in combination. (16) A system as in embodiment (15), wherein performance of a group of the computers configured as physical network devices forming a single logical device are power managed by reducing the performance and power consumption of each constituent physical device in predetermined equal increments or predetermined unequal increments. (17) A system as in embodiment (15), wherein network device loading and quality of service (QoS) are measured for a plurality of physical network devices organized as a single logical network device. (18) The computer system in embodiment (15), wherein the third processor clock frequency is less than the second processor clock frequency which is less than the first processor clock frequency. (19) The computer system in embodiment (18), wherein the second processor core voltage is less than the first processor core voltage. (20) The computer system in embodiment (19), wherein the third processor core voltage is less than the second processor core voltage. (21) The computer system in embodiment (15), wherein the third processor clock frequency is less than the second processor clock frequency which is less than the first processor clock frequency; and the second processor core voltage is less than the first processor core voltage. (22) The computer system in embodiment (2), wherein: each of the computers further being operable in (iii) a second mode having a second maximum performance level intermediate between the first maximum performance level and the third maximum performance level and a second power consumption rate intermediate between the first power consumption rate and the third power consumption rate; and each the computer including a local power manager determining an operating mode for itself selected from the first mode and the second mode based on processor internal activity information. (23) The computer system in embodiment (22), wherein the processor internal activity information comprising idle thread execution information. (24) The computer system in embodiment (22), wherein a transition from the first mode to the second mode is controlled locally within each the computer; and a transition from either the first mode or the second mode to the third mode are controlled globally by the power manager. (25) The computer system in embodiment (24), wherein a transition from the second mode to the first mode is controlled locally within each the computer; and a transition from the third mode to either the first mode or the second mode is controlled globally by the power manager. (26) The computer system in embodiment (15), wherein the third processor clock frequency is substantially zero or the third processor clock is turned off. (27) The computer system in embodiment (15), wherein the commands are generated and directed to the computers only when required to change an operating mode of the computers. (28) The computer system in any of embodiments (2 or 3), wherein the third mode is characterized by maintaining a processor core voltage to maintain processor state.

(29) A computer system comprising: a plurality of computers each having at least one processor and an activity monitor identifying a level of activity indicator for the at least one processor; each of the computers being operable in:

(i) a first mode having a first maximum performance level and a first power consumption rate, and (ii) a third mode having a third maximum performance level lower than the first maximum performance level and a third power consumption rate lower than the first power consumption rate; and a power manager: (i) coupled to each of the computers and receiving the level of activity information from each of the plurality of computers; (ii) analyzing the plurality of received level of activity information; (iii) determining an operating mode for each of the computers selected from the first mode and third mode based on the analyzed activity information and predetermined policies; and (iv) generating commands to each of the plurality of computers directing each of the plurality of computers to operate in the determined operating mode; each of the computers further being operable in (iii) a second mode having a second maximum performance level intermediate between the first maximum performance level and the third maximum performance level and a second power consumption rate intermediate between the first power consumption rate and the third power consumption rate; each the computer including a local power manager determining an operating mode for itself selected from the first mode and the second mode based on processor internal activity information; a transition from the first mode to the second mode is controlled locally within each the computer, and a transition from either the first mode or the second mode to the third mode are controlled globally by the power manager; and a transition from the second mode to the first mode is controlled locally within each the computer, and a transition from the third mode to either the first mode or the second mode is controlled globally by the power manager.

(30) A computer system comprising: a plurality of server computers each having at least one processor and an activity monitor identifying a level of activity for the at least one processor, the activity monitor comprising an activity monitor that monitors an activity selected from the set of activities consisting of: a program application layer activity, a network layer activity, a physical layer activity, and combinations thereof; each of the server computers being operable in: (i) a first mode having a first maximum performance level and a first power consumption rate, (ii) a second mode having a second maximum performance level lower than the first maximum performance level and a second power consumption rate lower than the first power consumption rate, and (iii) a third mode having a third maximum performance level lower than the second maximum performance level and a third power consumption rate lower than the second power consumption rate; and a power manager operative in a separate power manager computer: (i) coupled to each of the server computers and receiving the level of activity information from each of the plurality of computers; (ii) analyzing the plurality of received level of activity information; (iii) determining an operating mode for each of the server computers selected from the first mode, second mode, and third mode based on the analyzed activity information; and (iv) generating commands to each of the plurality of server computers directing each of the plurality of server computers to operate in the determined operating mode; the first mode operation is characterized by operating the processor at a first processor clock frequency and a first processor core voltage, the second mode operation is characterized by operating the processor at a second processor clock frequency and a second processor core voltage, and the third mode operation is characterized by operating the processor at a third processor clock frequency and a third processor core voltage; the second mode of operation being

further characterized in that the second processor clock frequency is lower than the first processor clock frequency and the second processor core voltage is equal to or less than the first processor core voltage so that in combination consuming less power than in the first mode, and the third processor clock frequency is lower than the second processor clock frequency and the third processor core voltage is no greater than the second processor core voltage so that in combination consuming less power than in the second mode; and a transition from the first mode to the second mode is controlled locally within each the computer; and a transition from either the first mode or the second mode to the third mode are controlled globally by the power manager.

(31) A method of operating computer system having a plurality of server computers, each server computer including at least one processor, and each computer being operable in a first mode having a first maximum performance level and a first power consumption rate, and a third mode having a third maximum performance level lower than the first maximum performance level and a third power consumption rate lower than the first power consumption rate; the method comprising: monitoring activity within the computers and identifying a level of activity for the at least one processor within the computers; analyzing the plurality of level of activity information; determining an operating mode for each of the computers selected from the first mode and third mode based on the analyzed activity information; and generating commands to each of the plurality of computers directing each of the plurality of computers to operate in the determined operating mode.

In a second group of innovations, the invention provides various embodiments associated with System and Method for Activity or Event Based Dynamic Energy Conserving Server Reconfiguration.

(1) An information processing system comprising: a frame or enclosure for mounting a plurality of devices; a backplane having a plurality of backplane electrical connectors disposed within the frame or enclosure; and a plurality of devices, each including a device electrical connector, matingly coupled to the backplane electrical connectors, the plurality of devices including at least one network device for coupling the system with an external network.

(2) A system as in embodiment (1), wherein the at least one network device comprises a device selected from the set of network devices consisting of a server device, a computer node device, a monitor node device, a management module, a server module, and combinations thereof. (3) A system as in embodiment (2), wherein the at least one network device includes a processor and a memory integral with or coupled to the processor. (4) A system as in embodiment (3), further comprising a network switch or network switching device. (5) A system as in embodiment (4), wherein the plurality of devices further comprises a device selected from the set of devices consisting of a power supply, a fan or fan module, and combinations thereof. (6) A system as in embodiment (1), wherein the at least one network device comprises at least one server computer having at least one processor and a power manager. (7) A system as in embodiment (6), wherein the power manager is integral with the server computer. (8) A system as in embodiment (6), wherein the power manager is separate from the server computer. (9) A system as in embodiment (6), wherein the or each server computer further comprises an activity monitor identifying a level of activity indicator for the at least one processor; and the or each server computer being operable in: (i) a first mode having a first maximum performance level and a first power consumption rate, and (ii) a third mode having a third

maximum performance level lower than the first maximum performance level and a third power consumption rate lower than the first power consumption rate; and the system further comprising: a power manager: (i) coupled to each of the computers and receiving the level of activity information from each of the plurality of computers; (ii) analyzing the plurality of received level of activity information; (iii) determining an operating mode for each of the computers selected from the first mode and third mode based on the analyzed activity information and predetermined policies; and (iv) generating commands to each of the plurality of computers directing each of the plurality of computers to operate in the determined operating mode. (10) A system as in embodiment (9), wherein: the or each server computer further being operable in: (iii) a second mode having a second maximum performance level intermediate between the first maximum performance level and the third maximum performance level and a second power consumption rate intermediate between the first power consumption rate and the third power consumption rate; and the power manager further determining an operating mode for each of the computers selected from the first mode, the second mode, and the third mode based on the analyzed activity information and the predetermined policies. (11) A system as in embodiment (1), wherein the system further comprises a power manager. (12) A system as in embodiment (10), wherein the system further comprises a switching module, and the power manager receives activity indicators for the switching module and controls an operating mode of the switching module in response thereto. (13) A system as in embodiment (10), wherein the computer comprises a server module that is power managed by adjusting processor performance to one or more of a predicted processor processing requirement and a measured processor processing requirement. (14) A system as in embodiment (13), wherein the predicted processor processing requirement is a Quality of Service (QoS) based requirement, and the measured processor processing requirement comprises a substantially real-time measured processor processing requirement. (15) A system as in embodiment (14), wherein the substantially real-time processor processing requirement comprises an idle thread execution detection and response thereto. (16) A system as in embodiment (10), wherein power (or energy) is conserved by controlling the computer based on a control procedure algorithm to enter a first level of power (energy) saving by adjusting the performance of the processor within the computer to substantially match the computer processor loading demand. (17) A system as in embodiment (10), wherein power (or energy) is conserved by controlling the plurality of computers in aggregate based on a control procedure algorithm and the policy to enter selected levels of power (energy) saving by adjusting the performance of the processors within the computers to one of the first mode, second mode, and third mode to substantially match the aggregate computer processor loading demands. (18) A system as in embodiment (10), wherein the power manager includes a control procedure algorithm implemented as software to implement a power on demand control procedure. (19) A system as in embodiment (10), wherein each computer is configurable as a particular type of network device. (20) A system as in embodiment (10), wherein the computer is configured as a network device selected from the set consisting of a web server, a streaming media server, a cache server, a file server, an application server, and a router. (21) A system as in embodiment (10), wherein at least selected ones of the computers are configurable as a combination type of network device, and wherein the network

device configured in the computer node is a network device selected from the set consisting of a web server, a streaming media server, a cache server, a file server, an application server, a router, and combinations thereof. (22) A system as in embodiment (21), wherein the network device is reconfigurable at any time based on types of activities detected within the network to which the network device is or may be connected. (23) A system as in embodiment (10), wherein at least one of the computers comprises a network device and the activity monitor for the network device comprises a network activity monitor that detects the types of activities present on a network to which the activity monitor is coupled. (24) A system as in embodiment (23), wherein the types of activities present on a network to which the activity monitor is coupled that are monitored by the activity monitor include volume of web pages served, volume of streaming media served, volume of files served, volume of applications served, volume of cached data served, amount of network traffic routed, and combinations thereof. (25) A system as in embodiment (22), wherein the reconfiguration of network device is initiated by any network device including the same network as is being reconfigured. (26) A system in embodiment (4) wherein a selected one of the plurality of computers designated as a master providing the power manager. (27) A system as in embodiment (10), wherein a selected one of the plurality of computers is designated as a master providing the power manager, and reconfiguration of a network device from one form of network device to another form of network device is initiated by any computer that has been designated as a master computer. (28) A system as in embodiment (27), wherein any computer may be designated as the master node. (29) A system as in embodiment (28), wherein a particular computer is designated as a master on the basis of its position within a chassis. (30) A system as in embodiment (28), wherein a particular computer node is designated as a master node on the basis of the order of power-up or boot completion. (31) A system as in embodiment (28), wherein reconfiguration of the computer comprises altering the software and/or firmware instructing the computer. (32) A system as in embodiment (28), wherein reconfiguration of the computer comprises altering the data organization of a data storage device integral with or coupled to the computer. (33) A system as in embodiment (28), wherein the data storage device comprises a hard disc drive based RAID storage array and altering the data organization comprises altering rad configuration of the data to provide better performance for the type of data being served. (34) A system as in embodiment (22), wherein the reconfiguration of a computer is initiated by a management module network device. (35) A system as in embodiment (10), wherein a plurality of computers of the same type are grouped together and treated as a single network device. (36) A system as in embodiment (35), wherein the group of network devices treated as a single network device is managed and controlled as a single network device. (37) A system as in embodiment (35), wherein the group of network devices treated as a single network device is power managed as a single network device. (38) A system as in embodiment (35), wherein the group of network devices treated as a single network device is monitored as a single network device. (39) A system as in embodiment (35), wherein the plurality of grouped network devices are electrically coupled via a backplane bus and the logical grouping of the plurality of network devices into a single logical network device is performed under control of software. (40) A system as in embodiment (39), wherein the software executes within a processor and memory associated within each network

device. (41) A system as in embodiment (35), wherein the plurality of network devices each comprise a server group. (42) A system as in embodiment (35), wherein the plurality of network devices each comprise a computer server module. (43) A system as in embodiment (42), wherein each computer server module is configured as a computer server module selected from the group consisting of a web server, a streaming media server, a cache server, a file server, an application server, a router, and combinations thereof. (44) A system as in embodiment (39), wherein the activity associated with each computer within a grouped logical network device may be monitored individually. (45) A system as in embodiment (39), wherein the network activity associated with all or any subset of physical network device within a grouped logical network device may be monitored as a composite or in aggregate. (46) A system as in embodiment (35), wherein grouping is accomplished by aggregating all of the activity in each computer and directing each computer in the logical group to operate at the same operating mode. (47) A system as in embodiment (10), wherein over a period of time the system will have sufficient over capacity that some of the computers will be directed to operate in the third mode, the policy taking into account the amount each of the computers have historically spent operating in at least one of the first, second, or third mode and selecting a computer to operate in the third mode based on historical data. (48) A system as in embodiment (47), wherein the computer selected to operate in the third mode is a computer that has the smallest cumulative duration operating in the third mode amongst the plurality of computers. (49) A system as in embodiment (47), wherein the computer selected to operate in the third mode is randomly selected from amongst the plurality of computers. (50) A system as in embodiment (47), wherein the computer selected to operate in the third mode is rotated sequentially amongst the plurality of computers. (51) A system as in embodiment (10), wherein the activity monitor comprises an activity monitor that monitors an activity selected from the set of activities consisting of: a program application layer activity, a network layer activity, a physical layer activity, and combinations thereof. (52) A system as in embodiment (51), wherein at the physical level the number of processor idle threads executed within a predetermined period of time are measured to determine processor loading and the processor performance is adjusted to by altering the operating mode to substantially match the level of processor loading. (53) A system as in embodiment (52), wherein the substantial matching of processor performance to processor loading is performed with a predetermined amount of additional processor performance beyond that needed to match the processor loading. (54) A system as in embodiment (53), wherein the predetermined amount of additional processor performance is between about one-percent and about five-percent additional performance. (55) The computer system in embodiment (10), wherein: the first mode operation is characterized by operating the processor at a first processor clock frequency and a first processor core voltage, the second mode operation is characterized by operating the processor at a second processor clock frequency and a second processor core voltage, and the third mode operation is characterized by operating the processor at a third processor clock frequency and a third processor core voltage; the second mode of operation being further characterized in that the second processor clock frequency and the second processor core voltage in combination consuming less power than the first processor clock frequency and the first processor core voltage in combination, and the third pro-

cessor clock frequency and the third processor core voltage in combination consuming less power than the second processor clock frequency and the second processor core voltage in combination. (56) A system as in embodiment (55), wherein performance of a group of the computers configured as physical network devices forming a single logical device are power managed by reducing the performance and power consumption of each constituent physical device in predetermined equal increments or predetermined unequal increments. (57) A system as in embodiment (56), wherein the unequal increments include placing one or more of the plurality of physical devices in the third mode operating mode. (58) A system as in embodiment (56), wherein the unequal increments include placing one or more of the plurality of physical devices in the second mode operating mode. (59) A system as in embodiment (56), wherein the unequal increments include placing one or more of the plurality of physical devices in a powered-off fourth mode.

(60) A system as in embodiment (56), wherein a composite performance of a logical network device is achieved by placing some physical network devices in the second mode and by placing others in a different mode. (61) The computer system in embodiment (8), wherein the activity monitor comprises a network layer activity monitoring TCP/IP protocol data packets; and processor performance is incrementally lowered by the power manager using the mode control until data packets start dropping indicating that the processor performance is at the limit of adequacy and then increasing the processor performance by a specified increment to act as a safety margin to provide reliable communication of the packets. (62) A system as in embodiment (61), wherein the specified increment is a one-percent to five percent increment. (63) A system as in embodiment (61), wherein the specified increment is a 0.1 percent to 10 percent increment. (64) The computer system in embodiment (9), wherein the activity monitor comprises an activity monitor that monitors an activity selected from the set of activities consisting of: a program application layer activity, a network layer activity, a physical layer activity, and combinations thereof. (65) The computer system in embodiment (64), wherein the application layer activity monitor comprises monitoring use of a port address within the computers, the monitoring including counting or measuring a number of times a specific port address is being requested within a predetermined period of time, and in response to that counting or measurement, placing a sufficient amount of computer performance to meet the performance requirement for each application requesting the port address. (66) A system as in embodiment (65), wherein the sufficient amount of network performance is provided by operating selected computer in a first predetermined performance having a predetermined power consumption and a second group of other selected physical network devices at a reduced second performance level having a power consumption lower than that of the first selected group. (67) A system as in embodiment (66), wherein the first predetermined performance is a maximum performance and the second predetermined performance is a second level power saving mode. (68) A system as in embodiment (66), wherein the first predetermined performance is a maximum performance and the second predetermined performance is a third level power saving mode. (69) A system as in embodiment (65), wherein the measurement is determined via a SNMP agent. (70) A system as in embodiment (9), wherein the power manager applies different policies for different application types including using different rules to determine and predict system performance

requirements. (71) A system as in embodiment (70), wherein the different application types comprise different server types. (72) A system as in embodiment (70), wherein the different rules comprise different measurement procedures. (73) A system as in embodiment (70), wherein the system performance requirements comprise processor performance requirements. (74) A system as in embodiment (70), wherein the system performance requirements comprise server loading performance requirements. (75) A system as in embodiment (70), wherein the application type comprises a network application. (76) A system as in embodiment (75), wherein the network application comprises a network file server (NFS) application. (77) The system in embodiment (76), wherein the computer comprises a network server, and a processor within the computer operates at a processor clock frequency just sufficient to maintain maximum rated communication over a predetermined network connection. (78) The system in embodiment (77), wherein the a predetermined network connection comprises a 100 Mbps ethernet connection. (79) A system as in embodiment (77), wherein the processor clock frequency is less than about 300 MHz. (80) A system as in embodiment (75), wherein the processor clock frequency is less than about 300 MHz. (81) The computer system in embodiment (10), wherein the activity indicator comprises a network quality of service indicator. (82) A system as in embodiment (10), wherein power is conserved by controlling each computer node to enter one of the second mode or the third mode using one or more of a quality of service based predictive processor performance reduction and a activity based measured performance requirement. (83) A system as in embodiment (82), wherein the activity based measured performance comprises an idle thread execution based activity measure. (84) A system as in embodiment (81), wherein a plurality of the computers are organized as a single logical network device, and network device loading and QoS are measured for logical network device. (85) A system as in embodiment (81), wherein within the single logical network device, at least some computers making up the logical network device enter the third mode while other of the physical network devices operate in one or more of the first and second modes. (86) A system as in embodiment (81), wherein the computers can enter a third mode directly or indirectly from either the first mode or the second mode. (87) A system as in embodiment (10), wherein when there is a requirement that one computer be placed in a lower power consumption mode, the computer selected for such lower power consumption is selected according to predetermined rules such that different computers are placed in lower power consumption mode each time such selection is required. (88) A system as in embodiment (87), wherein the predetermined rules provide for random selection of one of the computers. (89) A system as in embodiment (87), wherein the predetermined rules provide for cycling through the computers according to some predetermined ordering. (90) A system as in embodiment (89), wherein the predetermined rules provide for cycling through the computers according to some predetermined ordering in which computers having the lowest time in service are preferentially selected for continued operation and network devices having the longest time in service are selected for reduced power operation. (91) A system as in embodiment (90), wherein the reduced power operation includes being powered off. (92) A system as in embodiment (90), wherein the reduced power operation includes being placed in a suspend mode. (93) A system as in embodiment (10), wherein a computer placed in mode 3 is in a suspend state and may be woken up and placed in the first mode or

the second mode by any one of a plurality of events including by a wake on LAN signal event. (94) A system as in embodiment (10), wherein the transition from one power consumption mode to another power consumption mode is based on a procedure implemented in software. (95) A system as in embodiment (10), wherein the transition from one power consumption mode to another power consumption mode is based on a procedure implemented in hardware and software. (96) A system as in embodiment (10), wherein when there is need to operate fewer than all the computer, the particular computer or logical group of computers that is (are) turned off or placed in a reduced power consumption mode is cycled so that over time all of the network devices experience similar operating time histories. (97) A system as in embodiment (96), wherein the computers include a non-volatile memory for storing operational history. (98) A system as in embodiment (97), wherein the operational history includes a total operating time indicator. (99) A system as in embodiment (97), wherein the operational history includes a time in service indicator. (100) A system as in embodiment (97), wherein the operational history includes indicators for operational time at each operational mode. (101) A system as in embodiment (10), wherein at least some of the computers include a mass storage device including a rotatable storage device. (102) A system as in embodiment (101), wherein the rotatable mass storage device comprises a rotatable magnetic hard disk drive. (103) A system as in embodiment (101), wherein the rotatable mass storage device comprises a rotatable optical disk drive. (104) A system as in embodiment (101), wherein the rotatable mass storage device comprises a rotatable magneto-optical disk drive. (105) A system as in embodiment (101), wherein the rotatable mass storage device is power managed by controlling the rotation of a motor rotating the rotatable device, wherein the disc drive is not rotated when a computer associated with the drive is in a mode 3 operating condition. (106) A system as in embodiment (10), wherein the computers are configured as network server devices and a network load versus allocated network device performance profile is provided for each different type of network server device, and the performance level set for operation of the network device is established by reference to the profile. (107) A system as in embodiment (106), wherein the profile is implemented as an analytical expression executed in software or firmware. (108) A system as in embodiment (106), wherein the profile is implemented as a piecewise linear expression executed in software or firmware. (109) A system as in embodiment (106), wherein the profile is implemented as a look-up-table stored in a memory. (110) A system as in embodiment (10), wherein at least one of the computers comprises a network server device and the activity monitoring for the network server device comprises a monitoring of either the network device load or the network device quality of service (QoS); and wherein the monitoring is performed by the activity monitor or by a separate management computer, or both. (111) A system as in embodiment (10), wherein the system includes at least one temperature sensor within an enclosure holding the computers for monitoring and reporting the temperature proximate the sensor to a computers configured to monitor the temperature. (112) A system as in embodiment (10), wherein the system includes a plurality of temperature sensors within the enclosure reporting to one or more network devices. (113) A system as in embodiment (112), wherein the plurality of temperature sensors are spatially distributed to provide temperature monitoring of different network devices within the enclosure. (114) A system as in embodiment (112),

wherein the plurality of temperature sensors are spatially distributed to provide temperature monitoring of different network devices and power supplies within the enclosure. (115) A system as in embodiment (111), wherein when the temperature sensed by a temperature sensor is within a predetermined magnitude relationship of a first predetermined value at least one computer is transitioned to a lower power consumption state. (116) A system as in embodiment (115), wherein when the temperature sensed by a temperature sensor is within a predetermined magnitude relationship of a second predetermined value at least one computer is transitioned to a powered off state. (117) A system as in embodiment (111), wherein the operational mode of at least one computer is reduced to a lower power consuming and heat dissipating state in response to a temperature sensor reporting a temperature greater than or equal to a predetermined value. (118) A system as in embodiment (111), wherein after the power consumption operating mode has been lowered permitting the computer to be operated at a higher power consuming state when the temperature sensed is below a predetermined temperature value, the lower temperature value being selected to provide hysteresis and prevent oscillation between higher power state and lower powered state. (119) A system as in embodiment (115), wherein the lower power consumption state is achieved by lowering the clock frequency of the processor, the clock frequency of a bus coupling a processor to other components, or the operating voltage of the processor or other components. (120) A system as in embodiment (115), wherein the particular network device that is transitioned to a lower power consumption state is selected based on predetermined rules. (121) A system as in embodiment (120), wherein the predetermined rules include a quality of service indicator. (122) A system as in embodiment (121), wherein additional computer devices are sent to lower energy consuming modes if the temperature remains above a predetermined temperature value. (123) A system as in embodiment (10), wherein power consumption within the system is reduced by adjusting the number and motor speed of cooling fans responsible for cooling the computer. (124) A system as in embodiment (10), wherein a plurality of cooling fans are provided and operate under control of the power manager that controls each fan to provide cooling at the rate and location desired to maintain the computers within a predetermined temperature range. (125) A system as in embodiment (10), wherein the plurality of computers are disposed within a common enclosure and the system further comprising a plurality of temperature sensors and a plurality of cooling devices are also disposed within the enclosure, the plurality of temperature sensors communicating a temperature signal to a temperature control means and the control means adjusting the on/off status and operational parameters of the cooling units to extract heat according to predetermined rules. (126) A system as in embodiment (125), wherein the power manager comprises the temperature control means. (127) A system as in embodiment (125), wherein one of the computers within the enclosure comprises the temperature control means. (128) A system as in embodiment (9), wherein the system further includes a plurality of power supplies and the power supplies are controlled to maintain a required power output level and operate the power supplies at a preferred efficiency. (129) A system as in embodiment (128), wherein only selected ones of the plurality of power supplies are operated. (130) A system as in embodiment (128), wherein multiple ones of the power supplies are operated but each is operated at less than rated power output capacity. (131) A system as in

embodiment (10), wherein the temperature of the system is moderated by motor driven cooling fans and wherein a rotational speed of the motor drive cooling is adjusted to maintain a predetermined temperature range proximate a temperature sensor. (132) A system as in embodiment (10), wherein the rotational speed of a motor drive cooling is adjusted to maintain a predetermined temperature range within an enclosure.

(133) A power-conservative multi-node network device, comprising: an enclosure having a power supply and a back-plane bus; a plurality of hot-pluggable node devices in the form of printed circuit (PC) cards adapted for connection with the back-plane buss; and each the node device being reconfigurable in substantially real-time to adapt to changing conditions on the network.

(134) The network device in embodiment (133), wherein the plurality of hot-pluggable node devices comprise up to sixteen node devices. (135) The network device in embodiment (133), wherein each of the node devices includes power saving control features.

(136) A computer program product for use in conjunction with a computer system having a plurality of server computers, each server computer including at least one processor, and each computer being operable in a first mode having a first maximum performance level and a first power consumption rate, and a third mode having a third maximum performance level lower than the first maximum performance level and a third power consumption rate lower than the first power consumption rate, the computer program product comprising a computer readable storage medium and a computer program mechanism embedded therein, the computer program mechanism, comprising: a program module that directs at least one computer, to function in a specified manner, the program module including instructions for: monitoring activity within the computers and identifying a level of activity for the at least one processor within the computers; analyzing the plurality of level of activity information; determining an operating mode for each of the computers selected from the first mode and third mode based on the analyzed activity information; and generating commands to each of the plurality of computers directing each of the plurality of computers to operate in the determined operating mode.

(137) The computer program product of embodiment (136), wherein each of the computers further being operable in a second mode having a second maximum performance level intermediate between the first maximum performance level and the third maximum performance level and a second power consumption rate intermediate between the first power consumption rate and the third power consumption rate; and the determining an operating mode further comprising determining an operating mode for each of the computers selected from the first mode, the second mode, and the third mode based on the analyzed activity information. (138) The computer program product of embodiment (137), wherein a transition from the first mode to the second mode is controlled locally within each the computer; and a transition from either the first mode or the second mode to the third mode are controlled globally by the power manager. (139) The computer program product of embodiment (138), wherein a transition from the second mode to the first mode is controlled locally within each the computer; and a transition from the third mode to either the first mode or the second mode is controlled globally by the power manager.

In a third group of innovations, the invention provides various embodiments associated with System, Method,

Architecture, and Computer Program Product for Dynamic Power Management in a Computer System.

(1) In a computer system including at least one processing unit, a memory coupled to the at least one processing unit, and logic circuits coupled to the processing unit contributing to operation of the computer system, a method for controlling the operating mode and as a result the power consumption of the computer system between a plurality of operating modes each having a different electrical power consumption levels or ranges; the method comprising: while operating in a first selected operating mode exhibiting that first selected mode's characteristic power consumption range, (i) monitoring the computer system to detect the occurrence or non-occurrence of a first event; and (ii) transitioning the computer system from the first selected operating mode to a second selected operating mode exhibiting that second selected operating mode's power consumption range.

(2) The method in embodiment (1), wherein the first selected mode is a higher power consuming mode than the second selected mode. (3) The method in embodiment (1), wherein the first selected mode is a lower power consuming mode than the second selected mode. (4) The method in embodiment (1), wherein the computer system further comprises peripheral devices coupled to the at least one processing unit and the peripheral devices are power managed to reduce power consumption. (5) The method in embodiment (4), wherein the peripheral devices include a mass storage device storing data for retrieval of the data, and an output port for outputting selected portions of the stored data upon request. (6) The method in embodiment (1), wherein the first event comprises execution of a predetermined number of idle threads. (7) The method in embodiment (1), wherein the first event comprises execution of a single idle thread. (8) The method in embodiment (1), wherein the first event comprises execution of a predetermined plurality of idle threads. (9) The method in embodiment (1), wherein the first event comprises a wake on LAN signal event. (10) The method in embodiment (1), wherein the first event comprises the occurrence of some specified level of CPU processing capability availability that is derived from either an enumeration or a statistical evaluation of the idle thread or idle threads that are being or have been executed during some time period. (11) The method in embodiment (1), wherein one of the first and second events comprises a measured decrease in server load. (12) The method in embodiment (1), wherein one of the first and second events comprises a predicted decrease in server load. (13) The method in embodiment (1), wherein one of the first and second events comprises a measured decrease in processor tasking. (14) The method in embodiment (1), wherein one of the first and second events comprises a predicted decrease in processor tasking. (15) The method in embodiment (1), wherein one of the first and second events comprises a measured decrease in communication channel bandwidth. (16) The method in embodiment (1), wherein one of the first and second events comprises predicted decrease in communication channel bandwidth. (17) The method in embodiment (12), wherein the predicted decrease in server load is a prediction based at least in part on time of day. (18) The method in embodiment (12), wherein the predicted decrease in server load is a prediction based at least in part on a quality of service requirement. (19) The method in embodiment (12), wherein the predicted decrease in processor tasking is a prediction based at least in part on time of day. (20) The method in embodiment (12), wherein the predicted decrease in processor tasking is a prediction based at least in part type of content to be processed by the computer system.

(21) The method in embodiment (12), wherein the predicted decrease in server loading is a prediction based at least in part type of content to be served by the computer system. (22) The method in embodiment (12), wherein the manner of the prediction is further based on the content served by the server computer system. (23) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 1) in which the processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (24) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 2) in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (25) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 2') in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (26) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 2'') in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (27) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 2''') in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage just sufficient to maintain switching circuits in the processor unit at the processing unit clock frequency, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (28) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3) in which the processing unit is operated at a slow but non-zero frequency processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (29) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3') in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

(30) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3'') in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage. (31) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3''') in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than or equal to a maximum rated logic circuit operating voltage. (32) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3''''') in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than a maximum rated logic circuit operating voltage. (33) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 3''''''') in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and the logic circuit is operated at a substantially zero logic circuit clock frequency and at a logic circuit operating voltage that is just sufficient to maintain logic circuit operating state. (34) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 4) in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage. (35) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 4') in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and the logic circuit is powered off by removing the logic circuit clock and by removing the logic circuit operating voltage or by setting the logic circuit operating voltage below a level that will maintain state, except that a real-time clock and circuit for waking the logic circuit and the processing unit are maintained in operation. (36) The method in embodiment (1), wherein one of the first selected operating mode and the second selected operating mode comprises a mode (Mode 4'') in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and the logic circuit is powered off by removing the logic circuit clock and by removing the logic circuit operating voltage or by setting the logic circuit operating voltage below a level

that will maintain state, except that a circuit for waking the logic circuit and the processing unit are maintained in operation.

(37) The method in embodiment (1), further comprising: while operating in the second selected operating mode exhibiting that second selected mode's characteristic power consumption range, (i) monitoring the computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning the computer system from the second selected operating mode to a third selected operating mode exhibiting that third selected operating mode's power consumption range.

(38) The method in embodiment (1), wherein the first selected operating mode and the second selected operating mode comprises different operating modes selected from the set of operating modes consisting of: (i) a mode in which the processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; (ii) a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; and (iii) a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency.

(39) The method in embodiment (38), wherein the set further consists of a mode in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage.

(40) The method in embodiment (1), further comprising: while operating in the second selected operating mode exhibiting that second selected mode's characteristic power consumption range, (i) monitoring the computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning the computer system from the second selected operating mode to a third selected operating mode exhibiting that third selected operating mode's power consumption range.

(41) The method in embodiment (40), wherein the first selected operating mode and the second selected operating mode comprises different operating modes, and the second selected operating mode and the third selected operating mode comprise different operating modes, each of the first, second, and third operating modes being selected from the set of modes consisting of: (i) a mode in which the processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; (ii) a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; and (iii) a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor

unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency.

(42) The method in embodiment (41), wherein the set further consists of a mode in which the processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage. (43) A computer program product for use in conjunction with a computer system including at least one processing unit, a memory coupled to the at least one processing unit, and logic circuits coupled to the processing unit contributing to operation of the computer system, a method for controlling the operating mode and as a result the power consumption of the computer system between a plurality of operating modes each having a different electrical power consumption levels or ranges; the computer program product comprising a computer readable storage medium and a computer program mechanism embedded therein, the computer program mechanism, comprising: a program module that directs the computer system to function in a specified manner, the program module including instructions for: (i) monitoring the computer system to detect the occurrence or non-occurrence of a first event while operating in a first selected operating mode exhibiting that first selected mode's characteristic power consumption range; and (ii) transitioning the computer system from the first selected operating mode to a second selected operating mode exhibiting that second selected operating mode's power consumption range. (44) The computer program product in embodiment (43), wherein the program module further including instructions for: while operating in the second selected operating mode exhibiting that second selected mode's characteristic power consumption range, (i) monitoring the computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning the computer system from the second selected operating mode to a third selected operating mode exhibiting that third selected operating mode's power consumption range. (45) The computer program product in embodiment (44), wherein the first selected operating mode and the second selected operating mode comprises different operating modes, and the second selected operating mode and the third selected operating mode comprise different operating modes, each of the first, second, and third operating modes being selected from the set of modes consisting of: (i) a mode in which the processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; (ii) a mode in which the processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency; and (iii) a mode in which the processing unit is operated at a substantially zero frequency processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and the logic circuit is operated at substantially maximum rated logic circuit clock frequency. (46) The computer program product in embodiment (45), wherein the set further consists of a mode in which the processing unit is powered off by removing a processing unit clock frequency and a processing unit core voltage.

(47) A computer system comprising: at least one processing unit and a memory coupled to the at least one processing unit; and logic circuits coupled to the processing unit



contributing to operation of the computer system; a controller for controlling the operating mode and as a result, the power consumption of the computer system between a plurality of operating modes each having a different electrical power consumption levels or ranges; the controller 5 being operable while operating in a first selected operating mode exhibiting that first selected mode's characteristic power consumption range, (i) to monitor the computer system to detect the occurrence or non-occurrence of a first event; and (ii) to transition the computer system from the 10 first selected operating mode to a second selected operating mode exhibiting that second selected operating mode's power consumption range.

Those workers having ordinary skill in the art in light of the description provided will no doubt appreciate other aspects, features, and advantages of the inventive system, method, and software control. It will be appreciated that the 15 afore described procedures implemented in a computer environment may be implemented using hardware, software, and/or firmware, and combinations of these. The detection, analysis, monitoring, decision making, and control functions are particularly amenable to computer program software and 20 firmware implementations and may readily be implemented in a central processing unit (CPU), processor, controller, micro-controller, or other logic unit within or associated with the computers. Therefore the invention includes hardware and software implementations, and descriptions of 25 procedures and methods anticipate that such procedures and methods may be implemented as a computer program and computer program product.

The foregoing descriptions of specific embodiments of the present invention have been presented for purposes of illustration and description. They are not intended to be exhaustive or to limit the invention to the precise forms 30 disclosed, and obviously many modifications and variations are possible in light of the above teaching. The embodiments were chosen and described in order to best explain the principles of the invention and its practical application, to thereby enable others skilled in the art to best utilize the invention and various embodiments with various modifica- 35 tions as are suited to the particular use contemplated. It is intended that the scope of the invention be defined by the claims appended hereto and their equivalents. All publications, product or other data sheets, web-site content, and patent applications cited or referenced in this specification are herein incorporated by reference as if each individual 40 publication or patent application were specifically and individually indicated to be incorporated by reference.

I claim:

1. In a computer system comprising:

a plurality of computers operating as servers each computer having at least one processing unit, a memory coupled to said at least one processing unit, logic circuits coupled to said processing unit contributing to 55 operation of said computer system, and an activity monitor identifying a level of activity information for said at least one processor, a method for controlling the operating mode and as a result the power consumption each of said plurality of computers in said computer system between a plurality of operating modes each 60 having a different electrical power consumption levels or ranges; said method comprising:

configuring each of said computers to be operable in at least: (i) a first mode having a first maximum performance level and a first power consumption rate, (ii) a 65 third mode having a third maximum performance level lower than said first maximum performance level and a

third power consumption rate lower than said first power consumption rate, and (iii) a second mode having a second maximum performance level intermediate between said first maximum performance level and said 5 third maximum performance level and a second power consumption rate intermediate between said first power consumption rate and said third power consumption rate;

generating, by each of said plurality of computers, a level of activity information;

communicating the level of activity information from each of said plurality of computers to a power manager computer, the power manager computer being one of 10 the plurality of computers or a different computer;

receiving by the power manager computer said level of activity information from each of said plurality of 15 computers;

analyzing by said power manager computer said plurality of received level of activity information;

determining by said power manager computer an operating mode for each of said computers selected from said 20 first mode, said second mode, and third mode based on said analyzed activity information and predetermined policies;

generating commands to each of said plurality of computers directing each of said plurality of computers to 25 operate in said determined operating mode; and

while operating in a first selected operating mode exhibiting that first selected mode's characteristic power consumption range, (i) monitoring said computer system to detect the occurrence or non-occurrence of a first 30 event; and (ii) transitioning said computer system from said first selected operating mode to a second selected operating mode exhibiting that second selected operating mode's power consumption range.

2. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 2') in which said processing unit is operated at less than maximum rated processing unit 35 clock frequency and at less than or equal to a maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

3. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 2') in which said processing unit is operated at less than maximum rated processing unit 40 clock frequency and at less than a maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

4. The method in claim 1, wherein said activity monitor comprises an activity monitor that monitors an activity selected from the set of activities consisting of: a program application layer activity, a network layer activity, a physical layer activity, and combinations thereof; and 55

said application layer activity monitor comprises monitoring use of a port address within said computers, said monitoring including counting or measuring a number of times a specific port address is being requested within a predetermined period of time, and in response to that counting or measurement, placing a sufficient amount of computer performance to meet the performance requirement for each application requesting the port address.

5. The method of claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 2''') in which said processing unit is operated at less than maximum rated processing unit clock frequency and at less than a maximum rated processing unit core voltage just sufficient to maintain switching circuits in said processor unit at said processing unit clock frequency, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

6. The method in claim 1, wherein one of said first, selected operating mode and said second operating mode comprises a mode (Mode 3) in which said processing unit is operated at a slow but non-zero frequency processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

7. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 3') in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

8. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 3'') in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency and at a substantially maximum rated logic circuit operating voltage.

9. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 3''') in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and said logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than or equal to a maximum rated logic circuit operating voltage.

10. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a Mode (Mode 3''''') in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to maintain processor unit state, and said logic circuit is operated at a logic circuit clock frequency less than a maximum rated logic circuit clock frequency and at a logic circuit operating voltage that is less than a maximum rated logic circuit operating voltage.

11. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 3''''''') in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (processing unit clock stopped) and at a processing unit core voltage just sufficient to

maintain processor unit state, and said logic circuit is operated at a substantially zero logic circuit clock frequency and at a logic circuit operating voltage that is just sufficient to maintain logic circuit operating state.

12. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 4) in which said processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage.

13. The method in claim 1, wherein one of said first selected operating mode and second selected operating mode comprises a mode (Mode 4') in which said processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and said logic circuit is powered off by removing said logic circuit clock and by removing said logic circuit operating voltage or by setting said logic circuit operating voltage below a level that will maintain state, except that a real-time clock and circuit for waking said logic circuit and said processing unit are maintained in operation.

14. the method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 4'') in which said processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage; and said logic circuit is powered off by removing said logic circuit clock and by removing said logic circuit operating voltage or by setting said logic circuit operating voltage below a level that will maintain state, except that a circuit for waking said logic circuit and said processing unit are maintained in operation.

15. The method in claim 1, further comprising:

while operating in said second selected operating mode exhibiting that second selected mode's characteristic power consumption range, (i) monitoring said computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning said computer system from said second selected operating mode to a first selected operating mode exhibiting that third selected operating mode's power consumption range.

16. The method in claim 1, wherein said first selected operating mode and said second selected operating mode comprises different operating modes selected from the set of operating modes consisting of:

(i) a mode in which said processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency;

(ii) a mode in which said processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency; and

(iii) a mode in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency.

55

17. The method in claim 6, wherein said set further consists of a mode in which said processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage.

18. The method in claim 17 further comprising: while 5 operating in said second selected operating mode exhibiting that second selected mode's characteristic power consumption range, (i) monitoring said computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning said computer system from said second 10 selected operating mode to a third selected operating mode exhibiting that third selected operating mode's power consumption range.

19. The method in claim 18, wherein said first selected operating mode and said second selected operating mode 15 comprises different operating modes, and said second selected operating mode and said third selected operating mode comprise different operating modes, each of said first, second, and third operating modes being selected from the set of modes consisting of:

(i) a mode in which said processing unit is operating at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock 25 frequency;

(ii) a mode in which said processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage, and said logic circuit is 30 operated at substantially maximum rated logic circuit clock frequency; and

(iii) a mode in which said processing unit is operated at a substantially zero frequency processing unit clock frequency (clock stopped) and at less than or equal to a 35 maximum rated processing unit core voltage sufficient to maintain processor unit state, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency.

20. The method in claim 19, wherein said set further 40 consists of a mode in which said processing unit is powered off by removing a processing unit clock frequency (processing unit clock stopped) and a processing unit core voltage.

21. The method in claim 1, wherein the first selected mode is a higher power consuming mode than the second selected 45 mode.

22. The method in claim 1, wherein the first selected mode is a lower power consuming mode than the second selected mode.

23. The method in claim 1, wherein the computer system 50 further comprises peripheral devices coupled to said at least one processing unit and said peripheral devices are power managed to reduce power consumption.

24. The method in claim 23, wherein said peripheral devices include a mass storage device storing data for 55 retrieval of said data, and an output port for outputting selected portions of said stored data upon request.

25. The method in claim 1, wherein said first event comprises execution of a predetermined number of idle 60 threads.

26. The method in claim 1, wherein said first event comprises execution of a single idle thread.

27. The method in claim 1, wherein said first event comprises execution of a predetermined plurality of idle 65 threads.

28. The method in claim 1, wherein said first event comprises a wake on LAN single event.

56

29. The method in claim 1, wherein said first event comprises the occurrence of some specified level of CPU processing capability availability that is derived from either an enumeration or a statistical evaluation of the idle thread or idle threads that are being or have been executed during some time period.

30. The method in claim 1, wherein one of said first and second events comprises a measured decrease in server load.

31. The method in claim 1, wherein one of said first and second events comprises a predicted decrease in server load.

32. The method in claim 1, wherein one of said first and second events comprises a measured decrease in processor 10 tasking.

33. The method in claim 1, wherein one of said first and second events comprises a predicted decrease in processor 15 tasking.

34. The method in claim 1, wherein one of said first selected operating mode and said second selected operating mode comprises a mode (Mode 1) in which said processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock 20 frequency and at a substantially maximum rated logic circuit operating voltage.

35. A computer program product for use in conjunction with a computer system including a plurality of computers each having at least one processing unit, a memory coupled to said at least one processing unit, logic circuits coupled to 25 said processing unit contributing to operation of said computer system, and an activity monitor coupled with said processing unit a method for controlling the operating mode and as a result the power consumption of said computer system between a plurality of operating modes for each computer wherein each computer having a different electrical power consumption levels or ranges; the computer program product comprising a computer readable storage medium and a computer program mechanism embedded therein, the computer program mechanism, comprising:

a program module that directs said computer system to function in a specified manner, the program module including instructions for:

(i) monitoring said computer system including each of the plurality of computers to receive activity information from the activity monitors of each computer to detect the occurrence or non-occurrence of a first event while operating in a first selected operating mode exhibiting that first selected mode's characteristic power consumption range; and

(ii) transitioning said plurality of computers in said computer system on an individual computer basis based on said computer activity information from said first selected operating mode to a second selected operating mode exhibiting that second selected operating mode's power consumption range;

while operating in said second selected operating mode for any of said plurality of computers exhibiting that second selected mode's characteristic power consumption range, (i) monitoring each of said plurality of computers in said computer system to detect the occurrence or non-occurrence of a second event; and (ii) transitioning said plurality of computers within said computer system based on said activity information on a computer-by-computer basis from said second selected operating mode to a third selected operating mode exhibiting that third selected operating mode's power consumption range;

57

said first selected operating mode and said second selected operating mode comprises different operating modes, and said second selected operating mode and said third selected operating mode comprise different operating modes, each of said first, second, and third operating modes being selected from the set of modes consisting of:

- (i) a mode in which said processing unit is operated at substantially maximum rated processing unit clock frequency and at substantially maximum rated processing unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency;
- (ii) a mode in which said processing unit is operated at less than maximum rated processing unit clock frequency and at less than or equal to a maximum rated processing

58

unit core voltage, and said logic circuit is operated at substantially maximum rated logic circuit frequency; and

- (iii) a mode in which said processing unit is operated at a substantially zero frequency processing unit clock frequency and at less than or equal to a maximum rated processing unit core voltage sufficient to maintain processor unit state, and said logic circuit is operated at substantially maximum rated logic circuit clock frequency;

said set further comprises a mode in which said processing unit is powered off by removing a processing unit clock frequency and a processing unit core voltage.

\* \* \* \* \*

UNITED STATES PATENT AND TRADEMARK OFFICE  
**CERTIFICATE OF CORRECTION**

PATENT NO. : 7,228,441 B2  
APPLICATION NO. : 09/860210  
DATED : June 5, 2007  
INVENTOR(S) : Fung

Page 1 of 1

It is certified that error appears in the above-identified patent and that said Letters Patent is hereby corrected as shown below:

Claim 2, column 52, line 38; please delete “” after “Mode 2”;

Claim 6, column 53, line 11; please delete “,” after “first”;

Claim 35, column 57, line 1; please delete “ans” and insert --and--; line 11, please delete “curcuit” and insert --circuit--; line 14, please delete “oprated” and insert --operated--.

Signed and Sealed this

Seventh Day of August, 2007

A handwritten signature in black ink on a light gray dotted background. The signature reads "Jon W. Dudas" in a cursive style.

JON W. DUDAS

*Director of the United States Patent and Trademark Office*